

# Supplementary Materials for “Embedding learning”

This supplementary materials include proofs in “Embedding learning”.

**Proof of Lemma 1.** By the definition of sufficient embedding, we have

$$\begin{aligned}\min_f \mathbb{E}(L(f(\mathcal{X}(S)), Y)) &= \min_f \mathbb{E}(\mathbb{E}(L(f(\mathcal{X}(S)), Y)|S)) = \min_f \mathbb{E}(\mathbb{E}(L(f(\mathcal{X}(S)), Y)|\mathcal{X}(S))) \\ &= \min_g \mathbb{E}(\mathbb{E}(L(g(S), Y)|S)) = \min_g \mathbb{E}(L(g(S), Y)),\end{aligned}$$

where first equality follows from the law of conditioning, the second equality follows from Definition 1, and the second last equality holds by the pointwise minimization of both sides. The desired result then follows.  $\square$

**Proof of Lemma 2.** Note that  $\mathcal{X}^*$  is a sufficient embedding, then Lemma 1 yields that

$$\mathbb{E}(L(f^*(\mathcal{X}^*(S)), Y)) = \min_f \mathbb{E}(L(f(\mathcal{X}^*(S)), Y)) = \min_g \mathbb{E}(L(g(S), Y)),$$

Note further that we can choose  $g(s) = f(\mathcal{X}(s))$ , thus for any  $\mathcal{X}$  and  $f$ ,  $\min_g \mathbb{E}(L(g(S), Y)) \leq \mathbb{E}(L(f(\mathcal{X}(s)), Y))$ . Then for any  $C \geq U(\mathcal{X}^*)$ ,

$$\begin{aligned}\min_g \mathbb{E}(L(g(S), Y)) &\leq \min_{f; \mathcal{X}: U(\mathcal{X}) \leq C} \mathbb{E}(L(f(\mathcal{X}(s)), Y)) \\ &\leq \mathbb{E}(L(f^*(\mathcal{X}^*(S)), Y)) = \min_g \mathbb{E}(L(g(S), Y)),\end{aligned}$$

where the second last inequality follows from  $U(\mathcal{X}^*) \leq C$ . The desired result follows.  $\square$

**Proof of Theorem 1.** By Lemma 2,

$$\mathbb{P}(e(\hat{\boldsymbol{\theta}}) \geq c_1 \delta_n^{2\mu}) = \mathbb{P}\left(\mathbb{E}(L(\hat{f}(\hat{\mathcal{X}}(S)), Y) - L(f^*(\mathcal{X}^*(S)), Y)) \geq c_1 \delta_n^{2\mu}\right).$$

Then the desired result follows from Theorem 1 of [3]. This completes the proof.  $\square$

**Proof of Corollary 1.** Without loss of generality, assume that the Lipschitz constant is 1 for  $V(z)$ , then it suffices to verify Assumption A3. Let  $\mathcal{V}(r) = \{V(yf(\mathcal{X}(s))) - V(y\bar{f}(\bar{\mathcal{X}}(s))) : f \in \mathcal{F}, J(f) \leq \bar{J}r, U(\mathcal{X}) \leq C\}$  and  $\mathcal{F}(r) = \{\boldsymbol{\beta}^\top \mathbf{x} : \|\boldsymbol{\beta}\|_2^2 \leq \bar{J}r\}$ . Then we compute the cover number  $\mathcal{N}(\mathcal{V}(r), \nu)$  of the functional space. For  $\boldsymbol{\beta}$  and  $\tilde{\boldsymbol{\beta}}$  in  $\mathcal{F}(r)$ ,  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$  in  $\mathcal{D}_C$ ,

$$\begin{aligned} & \sup_{S, Y} |V(Y\tilde{\boldsymbol{\beta}}^\top \tilde{\mathcal{X}}(S)) - V(Y\boldsymbol{\beta}^\top \mathcal{X}(S))| \leq \max_{u=1, \dots, |S|} |\tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_u - \boldsymbol{\beta}^\top \tilde{\mathbf{x}}_u + \boldsymbol{\beta}^\top \tilde{\mathbf{x}}_u - \boldsymbol{\beta}^\top \mathbf{x}_u| \\ & \leq \max_{u=1, \dots, |S|} \|\tilde{\mathbf{x}}_u\|_2 \|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2 + \|\boldsymbol{\beta}\|_2 \max_{u=1, \dots, |S|} \|\tilde{\mathbf{x}}_u - \mathbf{x}_u\|_2 \\ & \leq B\sqrt{p} \|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_2 + \sqrt{\bar{J}pr} \|\text{vec}(\tilde{\mathbf{X}}) - \text{vec}(\mathbf{X})\|_\infty, \end{aligned}$$

where  $\mathbf{x}_u = \mathcal{X}(s_u)$  is the  $u$ -th column vector of  $\mathbf{X}$ , the second inequality follows from the Cauchy-Schwarz inequality. Then, it suffices to compute the covering number of  $\mathcal{F}(r)$  and  $\mathcal{X}_C$  separately, i.e.,

$$\begin{aligned} H_\infty(\nu, \mathcal{V}(r)) & \leq \log \mathcal{N}(\nu/(2B\sqrt{p}), \mathcal{F}(r)) + \log \mathcal{N}(\nu/(2\sqrt{\bar{J}pr}), \mathcal{X}_C) \\ & \leq p \left( \log \left( \frac{6B\sqrt{\bar{J}pr}}{\nu} \right) + \sum_{u=1}^{|S|} \max \left( \log \left( \frac{\sqrt{2\bar{J}pr} \|\mathbf{P}\|_\infty \min(\sqrt{C\sigma_u^{-1}}, B\|\mathbf{P}\|_1)}{\nu} \right), 0 \right) \right), \end{aligned}$$

where  $\mathcal{X}_C = \{\text{vec}(\mathbf{X}) \in \mathbb{R}^{p|S|} : \text{Tr}(\mathbf{X}\mathbf{Q}\mathbf{X}^T) \leq 2C, \|\text{vec}(\mathbf{X})\|_\infty \leq B\}$ , and the last inequality from Lemma 7. To apply Theorem 1, we verify the required entropy condition for  $\mathcal{V}(r)$ ,  $\sup_{r \geq 2} \phi(\epsilon_n, r) \leq c_3 n^{1/2}$ , where  $\phi(\epsilon_n, r) = \int_{c_5 \zeta}^{c_4^{1/2} \zeta^{\kappa/2}} H_\infty^{1/2}(u, \mathcal{V}(r)) du / \zeta$ ,  $\zeta = \min(\epsilon_n^2 + \lambda \bar{J}(r/2 - 1), 1)$ . Note that  $\phi(\epsilon_n, r)$  is non-increasing in  $r \geq 2$  for any fixed  $\epsilon_n > 0$ . Then

$\sup_{r \geq 2} \phi(\epsilon_n, r) = \phi(\epsilon_n, 2)$  is bounded by

$$\begin{aligned}
& p^{1/2} \left( \log \left( \frac{6B\sqrt{2\bar{J}p}}{\epsilon_n^2} \right) + \sum_{u=1}^{|S|} \max \left( \log \left( \frac{2\sqrt{\bar{J}p} \|\mathbf{P}\|_\infty \min(\sqrt{C\sigma_u^{-1}}, B\|\mathbf{P}\|_1)}{\epsilon_n^2} \right), 0 \right) \right)^{1/2} \epsilon_n^{\kappa-2} \\
& \leq p^{1/2} \left( \log \left( \frac{6B\sqrt{2\bar{J}p}}{\epsilon_n^2} \right) + \sum_{u=1}^{K^*} \log \left( \frac{2\sqrt{\bar{J}p} B \|\mathbf{P}\|_1 \|\mathbf{P}\|_\infty}{\epsilon_n^2} \right) \right)^{1/2} \epsilon_n^{\kappa-2} \\
& \leq p^{1/2} (1 + K^*)^{1/2} \left( \log \left( \frac{6B\sqrt{2\bar{J}p} (1 + \|\mathbf{P}\|_1 \|\mathbf{P}\|_\infty)}{\epsilon_n^2} \right) \right)^{1/2} \epsilon_n^{\kappa-2} \leq c_3 n^{1/2}, \tag{1}
\end{aligned}$$

where the second inequality follows from the definition of  $K^*$ , and the last inequality implies from the fact that  $\epsilon_n^2 \sim \left( \frac{(K^*+1)p^*}{n} \log \left( \frac{nB\sqrt{\bar{J}p^*}(1+\|\mathbf{P}\|_1\|\mathbf{P}\|_\infty)}{(K^*+1)p^*} \right) \right)^{\frac{1}{2-\kappa}}$ .  $\square$

**Proof of Corollary 2.** Without loss of generality, we assume  $V(z)$  is Lipschitz continuous with the Lipschitz constant equal to 1. We bound the approximation error and estimation error separately.

For the estimation error, let  $\mathcal{I} = [-B, B]^{p^*}$ . For any  $f_M, \tilde{f}_M \in \mathcal{F}(r) = \{f \in \mathcal{F}, J(f) \leq r\bar{J}\}$ , with  $\mathcal{F}$  being a space of neural networks and  $\mathbf{X}, \tilde{\mathbf{X}} \in \mathcal{D}_C$ ,

$$\begin{aligned}
& \sup_{S, Y} |V(Y\tilde{f}_M(\tilde{\mathcal{X}}(S))) - V(Yf_M(\mathcal{X}(S)))| \\
& \leq \max_{u=1, \dots, |S|} |\tilde{f}_M(\tilde{\mathbf{x}}_u) - f_M(\mathbf{x}_u)| \leq \max_{u=1, \dots, |S|} (|\tilde{f}_M(\tilde{\mathbf{x}}_u) - f_M(\tilde{\mathbf{x}}_u)| + |f_M(\tilde{\mathbf{x}}_u) - f_M(\mathbf{x}_u)|) \\
& \leq \sup_{\mathbf{x} \in \mathcal{I}} |\tilde{f}_M(\mathbf{x}) - f_M(\mathbf{x})| + \max_{u=1, \dots, |S|} |f_M(\tilde{\mathbf{x}}_u) - f_M(\mathbf{x}_u)| \\
& \leq \sup_{\mathbf{x} \in \mathcal{I}} |\tilde{f}_M(\mathbf{x}) - f_M(\mathbf{x})| + \|\mathbf{A}_M\|_{1,1} \cdots \|\mathbf{A}_1\|_{1,1} \max_{u=1, \dots, |S|} \|\tilde{\mathbf{x}}_u - \mathbf{x}_u\|_\infty \\
& \leq \sup_{\mathbf{x} \in \mathcal{I}} |\tilde{f}_M(\mathbf{x}) - f_M(\mathbf{x})| + (r\bar{J}/M)^M \|\text{vec}(\tilde{\mathbf{X}}) - \text{vec}(\mathbf{X})\|_\infty,
\end{aligned}$$

where the last inequality follows from the means inequality, and the second last inequality

follows from the fact that

$$\begin{aligned}
& \max_{1 \leq j \leq h_m} |(\mathbf{f}_m(\tilde{\mathbf{x}}_u) - \mathbf{f}_m(\mathbf{x}_u))_j| = \max_{1 \leq j \leq h_m} |\sigma_j(\mathbf{A}_m \mathbf{f}_{m-1}(\tilde{\mathbf{x}}_u) + \mathbf{b}_m) - \sigma_j(\mathbf{A}_m \mathbf{f}_{m-1}(\mathbf{x}_u) + \mathbf{b}_m)| \\
& \leq \max_{1 \leq j \leq h_m} |(\mathbf{A}_m)_j \mathbf{f}_{m-1}(\tilde{\mathbf{x}}_u) - (\mathbf{A}_m)_j \mathbf{f}_{m-1}(\mathbf{x}_u)| \\
& \leq \max_{1 \leq j \leq h_m} \|(\mathbf{A}_m)_j\|_1 \max_{1 \leq j \leq h_{m-1}} |(\mathbf{f}_{m-1}(\tilde{\mathbf{x}}_u) - \mathbf{f}_{m-1}(\mathbf{x}_u))_j| \\
& \leq \|\mathbf{A}_m\|_{1,1} \max_{1 \leq j \leq h_{m-1}} |(\mathbf{f}_{m-1}(\tilde{\mathbf{x}}_u) - \mathbf{f}_{m-1}(\mathbf{x}_u))_j|.
\end{aligned}$$

Then it suffices to bound the entropy of  $\mathcal{F}(r)$  and  $\mathcal{D}_C$ . Since  $\max_{s \in \mathcal{S}} \|\mathcal{X}(s)\|_\infty \leq B$ ,

$$\begin{aligned}
H_\infty(\nu, \mathcal{V}(r)) & \leq \log \mathcal{N}(\nu/2, \mathcal{F}(r)) + \log \mathcal{N}\left(\frac{\nu}{2(r\bar{J}/M)^M}, \mathcal{D}_C\right) \\
& \leq |\Theta_N| \log \left( \frac{6r\bar{J}}{\nu} \left( MB \left( \frac{2r\bar{J}}{M-1} \right)^{M-1} + \left( \frac{r\bar{J}(1 - (r\bar{J})^M)}{1 - r\bar{J}} \right)^2 \right) \right. \\
& \quad \left. + p^* \sum_{u=1}^{|S|} \max \left( \log \left( \frac{\sqrt{2}\|\mathbf{P}\|_\infty (r\bar{J}/M)^M \min(\sqrt{C\sigma_u^{-1}}, B\|\mathbf{P}\|_1)}{\nu} \right), 0 \right) \right) \\
& \leq |\Theta_N| \log \left( \frac{6(MB+1)(r\bar{J})^{2M+1}}{\nu} \right. \\
& \quad \left. + p^* \sum_{u=1}^{|S|} \max \left( \log \left( \frac{\sqrt{2}\|\mathbf{P}\|_\infty (r\bar{J}/M)^M \min(\sqrt{C\sigma_u^{-1}}, B\|\mathbf{P}\|_1)}{\nu} \right), 0 \right) \right),
\end{aligned}$$

where the last inequality follows from Lemma 6. To apply Theorem 1, we verify the required entropy condition for  $\mathcal{V}(r)$ ,  $\sup_{r \geq 2} \phi(\epsilon_n, r) \leq c_3 n^{1/2}$ . Note that  $\phi(\epsilon_n, r)$  is non-increasing in  $r$  for  $r \geq 2$  for any fixed  $\epsilon_n > 0$ ,

$$\begin{aligned}
\sup_{r \geq 2} \phi(\epsilon_n, 2) & = \phi(\epsilon_n, 2) \leq \left( |\Theta_N| \log \left( \frac{6MB(2\bar{J})^{2M+1}}{\epsilon_n^2} \right) \right. \\
& \quad \left. + p^* \sum_{u=1}^{|S|} \max \left( \log \left( \frac{\sqrt{2}\|\mathbf{P}\|_\infty (r\bar{J}/M)^M \min(\sqrt{C\sigma_u^{-1}}, B\|\mathbf{P}\|_1)}{\epsilon_n^2} \right), 0 \right) \right)^{1/2} \epsilon_n^{\kappa-2} \\
& \leq \left( |\Theta_N| \log \left( \frac{6MB(2\bar{J})^{2M+1}}{\epsilon_n^2} \right) + p^* K^* \log \left( \frac{\sqrt{2}B(2\bar{J}/M)^M \|\mathbf{P}\|_1 \|\mathbf{P}\|_\infty}{\epsilon_n^2} \right) \right)^{1/2} \epsilon_n^{\kappa-2} \\
& \leq (|\Theta_N| + p^* K^*)^{1/2} \left( \log \left( \frac{6B(2\bar{J})^{2M+1}(M + \|\mathbf{P}\|_1 \|\mathbf{P}\|_\infty)}{\epsilon_n^2} \right) \right)^{1/2} \leq c_3 n^{1/2},
\end{aligned}$$

where  $\epsilon_n^2 \sim \left( \frac{|\Theta_N| + K^* p^*}{n} \log \left( \frac{nB(2\bar{J})^{2M}(M + \|\mathbf{P}\|_1 \|\mathbf{P}\|_\infty)}{M(|\Theta_N| + K^* p)} \right) \right)^{\frac{1}{2-\kappa}}$ , the second inequality follows from the fact that  $2\sqrt{C}\|\mathbf{P}\|_\infty(2\bar{J}_M/M)^M/\sqrt{\sigma_k} \leq \epsilon_n^2$ , for any  $k > K^*$ .

Next, we bound the approximation error  $v_n^2 = e_V(\bar{\boldsymbol{\theta}})$  based on Theorem 1 of [6]. First, we transform  $\mathcal{X}^*(s)$  from  $[-B^*, B^*]^{p^*}$  to  $[0, 1]^{p^*}$ . To proceed, let  $\mathcal{X}_{B^*}^*(s) = \frac{\mathcal{X}^*(s) + B^*}{2B^*\tau} \in [0, 1/\tau]^{p^*} \subset [0, 1]^{p^*}$  and  $\mathcal{X}_{B^*}^* = \{\mathcal{X}_{B^*}^*(s)\}_{s \in \mathcal{S}}$ . Since  $f^* \in \mathcal{W}_\infty^\alpha([-B^*, B^*]^{p^*})$ , then  $f_{B^*}^*(\mathbf{x}) = f^*(2\tau B^* \mathbf{x} - B^*) \in \mathcal{W}_\infty^\alpha([0, 1/\tau]^{p^*})$ , where  $\tau = \max(1, \Gamma/U(\mathcal{X}^*))^{1/2} \geq 1$ , and  $\Gamma = p|\mathcal{S}|\text{Tr}(\mathbf{Q})$ . Then the desired result follows from the fact that

$$\begin{aligned} v_n^2 = e_V(\bar{\boldsymbol{\theta}}) &= \mathbb{E} \left( V(Y \bar{f}_M(\bar{\mathcal{X}}(S))) - V(Y f^*(\mathcal{X}^*(S))) \right) \\ &= \mathbb{E} \left( V(Y \bar{f}_M(\mathcal{X}_{B^*}^*(S))) - V(Y f_{B^*}^*(\mathcal{X}_{B^*}^*(S))) \right) \leq \|\bar{f}_M - f_{B^*}^*\|_\infty \\ &\leq a_1 \max \left( \frac{1}{e^{M-1}}, |\Theta_N|^{-\frac{\alpha}{p^*}} \log \left( \frac{n}{M(|\Theta_N| + K^* p^*)} \right) \right), \end{aligned}$$

where  $a_1 > 0$  is a constant, the second inequality follows from Theorem 1 of [6], and the first equality holds by taking  $\bar{\mathcal{X}} = \mathcal{X}_{B^*}^*$ , and  $\mathcal{X}_{B^*}^* \in \mathcal{D}_C$ , since

$$\begin{aligned} U(\mathcal{X}_{B^*}^*) &\leq \frac{1}{\max(1, \Gamma/U(\mathcal{X}^*))} \sup_{\mathbf{X} \in [0, 1]^{|\mathcal{S}| \times |\mathcal{S}|}} \text{Tr}(\mathbf{X} \mathbf{Q} \mathbf{X}^\top) \\ &\leq \sup_{\mathbf{X} \in [0, 1]^{|\mathcal{S}| \times |\mathcal{S}|}} \frac{1}{\max(1, \Gamma/U(\mathcal{X}^*))} \sum_i \sigma_i \sigma_i(\mathbf{X}^\top \mathbf{X}) \\ &\leq \frac{p^* |\mathcal{S}| \text{Tr}(\mathbf{Q})}{\max(1, \Gamma/U(\mathcal{X}^*))} = \min(\Gamma, U(\mathcal{X}^*)) \leq C, \end{aligned}$$

where  $\sigma_i(\mathbf{X}^\top \mathbf{X})$  is the  $i$ -th eigenvalue of  $\mathbf{X}^\top \mathbf{X}$ . This completes the proof.  $\square$

**Proof of Lemma 3.** Note that (14) and (17) are convex minimization. Then convergence of Algorithm 1 follows from Theorem 4.7 of [1]. Moreover, local minimum of the solution follows from [5]. This completes the proof.  $\square$

**Proof of Lemma 4.** Note that the definition of a stationary point in Proposition 4 of [2] is in fact a local minimizer. By Proposition 4 there, we obtain the desired result. This completes the proof.  $\square$

**Proof of Lemma 5.** Let  $\mathbf{g} = (\mathbb{E}(Y|S = s_1), \dots, \mathbb{E}(Y|S = s_{|S|}))^\top$  and  $\mathcal{B}_{|S|}$  is a unit  $L_2$ -ball in  $\mathbb{R}^{|S|}$ . Then, (27) is equivalent to the existence of linear equations  $\mathbf{g} = \widetilde{\mathbf{X}}^\top \boldsymbol{\beta}$  for any  $\mathbf{g} \in \mathcal{B}_{|S|}$  and a given embedding matrix  $\widetilde{\mathbf{X}}$ . By the Rouché-Capelli Theorem [4], it is equivalent to  $\text{Rank}(\widetilde{\mathbf{X}}) = \text{Rank}([\widetilde{\mathbf{X}}; \mathbf{g}])$  for any  $\mathbf{g} \in \mathcal{B}_{|S|}$ , where  $[\widetilde{\mathbf{X}}; \mathbf{g}]$  is the augmented matrix. Then, the sufficiency is established. Now, we prove the necessity by contradiction. Suppose  $\text{Rank}(\widetilde{\mathbf{X}}) < |S|$ . Then,  $\text{Span}(\widetilde{\mathbf{X}})^\perp \neq \emptyset$ , and thus  $\text{Rank}([\widetilde{\mathbf{X}}; \mathbf{g}]) > \text{Rank}(\widetilde{\mathbf{X}})$  when  $\mathbf{g} = \mathbf{z}/\|\mathbf{z}\|_2$ ,  $\mathbf{z} \in \text{Span}(\widetilde{\mathbf{X}})^\perp$ . This contradicts to that (27). This completes the proof.  $\square$

**Lemma 6.** Consider a function space of  $M$ -depth ReLU neural networks defined as

$$\mathcal{F}(r) = \{\mathbf{f}_M(\mathbf{x}), \text{ with } (\mathbf{f}_m(\mathbf{x}) = \boldsymbol{\sigma}(\mathbf{A}_m \mathbf{f}_{m-1}(\mathbf{x}) + \mathbf{b}_m))_{m=1}^M : (\mathbf{A}, \mathbf{b}) \in \boldsymbol{\Theta}_N(r)\},$$

where  $\boldsymbol{\sigma}(\mathbf{z}_m) = (\sigma(z_{m,1}), \dots, \sigma(z_{m,h_m}))^\top$ ,  $\sigma(\cdot)$  is a rectified linear unit (ReLU),  $\boldsymbol{\Theta}_N(r) = \{(\mathbf{A}, \mathbf{b}) : \mathbf{A} = (\mathbf{A}_m)_{m=1}^M, \mathbf{b} = (\mathbf{b}_m)_{m=1}^M, \|\mathbf{A}\|_{1,1} + \|\mathbf{b}\|_{1,1} = \sum_{m=1}^M (\sum_{j=1}^{h_m} \|\mathbf{A}_{m,j}\|_1 + \|\mathbf{b}_m\|_1) \leq r\}$ , and  $\mathbf{x} \in \mathcal{I} \subset [-B, B]^p$ , then for any  $\tilde{\mathbf{f}}_M$  and  $\mathbf{f}_M$  in  $\mathcal{F}(r)$ , we have

$$H_\infty(\mathcal{F}(r), \nu) \leq |\boldsymbol{\Theta}_N| \log \left( \frac{6r}{\nu} \left( MB \left( \frac{2r}{M-1} \right)^{M-1} + \left( \frac{r(1-r^M)}{1-r} \right)^2 \right) \right),$$

where  $|\boldsymbol{\Theta}_N|$  is the number of parameters in the neural network.

**Proof of Lemma 6.** Let  $\Lambda_m = \|\mathbf{A}_m\|_{1,1}$  and  $\tilde{\Lambda}_m = \|\tilde{\mathbf{A}}_m\|_{1,1}$ , for any  $\mathbf{f}_M$  and  $\tilde{\mathbf{f}}_M$  in  $\mathcal{F}(r)$ , we have

$$\begin{aligned} \max_{1 \leq j \leq h_m} |(\mathbf{f}_m(\mathbf{x}) - \tilde{\mathbf{f}}_m(\mathbf{x}))_j| &= \max_{1 \leq j \leq h_m} |\boldsymbol{\sigma}_j(\mathbf{A}_m^\top \mathbf{f}_{m-1}(\mathbf{x}) + \mathbf{b}_m) - \boldsymbol{\sigma}_j(\tilde{\mathbf{A}}_m^\top \tilde{\mathbf{f}}_{m-1}(\mathbf{x}) + \tilde{\mathbf{b}}_m)| \\ &\leq \max_{1 \leq j \leq h_m} |(\mathbf{A}_m)_j^\top \mathbf{f}_{m-1}(\mathbf{x}) - (\tilde{\mathbf{A}}_m)_j^\top \tilde{\mathbf{f}}_{m-1}(\mathbf{x})| + \|\mathbf{b}_m - \tilde{\mathbf{b}}_m\|_\infty \\ &\leq \|\mathbf{A}_m - \tilde{\mathbf{A}}_m\|_{1,1} \max_{1 \leq j \leq h_m} |(\mathbf{f}_{m-1}(\mathbf{x}))_j| + \tilde{\Lambda}_m \max_{1 \leq j \leq h_m} |(\mathbf{f}_{m-1}(\mathbf{x}) - \tilde{\mathbf{f}}_{m-1}(\mathbf{x}))_j| + \|\mathbf{b}_m - \tilde{\mathbf{b}}_m\|_\infty \\ &\leq \|\mathbf{A}_m - \tilde{\mathbf{A}}_m\|_{1,1} (\Lambda_{m-1} \cdots \Lambda_1 \|\mathbf{x}\|_\infty + \frac{r(1-r^m)}{1-r}) \\ &\quad + \tilde{\Lambda}_m \max_{1 \leq j \leq h_m} |(\mathbf{f}_{m-1}(\mathbf{x}) - \tilde{\mathbf{f}}_{m-1}(\mathbf{x}))_j| + \|\mathbf{b}_m - \tilde{\mathbf{b}}_m\|_\infty, \end{aligned}$$

where the second last inequality follows from the Cauchy-Schwarz inequality, and the last inequality follows from the fact that

$$\begin{aligned} \max_{1 \leq j \leq h_m} |(\mathbf{f}_m(\mathbf{x}))_j| &\leq \max_{1 \leq j \leq h_m} |\boldsymbol{\sigma}_j(\mathbf{A}_m^\top \mathbf{f}_{m-1}(\mathbf{x}) + \mathbf{b}_m)| \leq \max_{1 \leq j \leq h_m} |(\mathbf{A}_m^\top \mathbf{f}_{m-1}(\mathbf{x}) + \mathbf{b}_m)_j| \\ &\leq \Lambda_m \max_{1 \leq j \leq h_m} |\mathbf{f}_{m-1}(\mathbf{x})| + \|\mathbf{b}_m\|_\infty. \end{aligned}$$

Next, multiplying both sides by  $\tilde{\Lambda}_M \cdots \tilde{\Lambda}_{m+1}$  and summing up from  $m = 1$  to  $m = M$  yield that

$$\begin{aligned} \max_{1 \leq j \leq h_m} |f_M(\mathbf{x}) - \tilde{f}_M(\mathbf{x})| &\leq (\Lambda_1 \cdots \Lambda_{M-1} + \Lambda_1 \cdots \Lambda_{M-2} \tilde{\Lambda}_M \cdots + \tilde{\Lambda}_2 \cdots \tilde{\Lambda}_M) \|\mathbf{A} - \tilde{\mathbf{A}}\|_{1,1} \|\mathbf{x}\|_\infty \\ &\quad + (1 + \tilde{\Lambda}_M \cdots + \tilde{\Lambda}_2 \cdots \tilde{\Lambda}_M) \left( \frac{r(1-r^M)}{1-r} \|\mathbf{A} - \tilde{\mathbf{A}}\|_{1,1} + \|\mathbf{b} - \tilde{\mathbf{b}}\|_\infty \right) \\ &\leq \left( \left( \frac{\Lambda_1 + \cdots + \Lambda_M}{M-1} \right)^{M-1} \cdots + \left( \frac{\tilde{\Lambda}_2 + \cdots + \tilde{\Lambda}_M}{M-1} \right)^{M-1} \right) \|\mathbf{A} - \tilde{\mathbf{A}}\|_{1,1} \|\mathbf{x}\|_\infty \\ &\quad + \left( 1 + r + \cdots + \left( \frac{r}{M-1} \right)^{M-1} \right) \left( \frac{r(1-r^M)}{1-r} \|\mathbf{A} - \tilde{\mathbf{A}}\|_{1,1} + \|\mathbf{b} - \tilde{\mathbf{b}}\|_\infty \right) \\ &\leq \left( M \left( \frac{2r}{M-1} \right)^{M-1} \|\mathbf{x}\|_\infty + \left( \frac{r(1-r^M)}{1-r} \right)^2 \right) \left( \|\mathbf{A} - \tilde{\mathbf{A}}\|_{1,1} + \|\mathbf{b} - \tilde{\mathbf{b}}\|_{1,1} \right). \end{aligned}$$

Taking  $\sup_{\mathbf{x} \in \mathcal{I}}$  on both sides yields that

$$\begin{aligned} H_\infty(\mathcal{F}(r), \nu) &\leq \log \mathcal{N} \left( \boldsymbol{\Theta}_N(r), \nu / \left( 2MB \left( \frac{2r}{M-1} \right)^{M-1} + 2 \left( \frac{r(1-r^M)}{1-r} \right)^2 \right) \right) \\ &\leq |\boldsymbol{\Theta}_N| \log \left( \frac{6r}{\nu} \left( MB \left( \frac{2r}{M-1} \right)^{M-1} + \left( \frac{r(1-r^M)}{1-r} \right)^2 \right) \right). \end{aligned}$$

This completes the proof.  $\square$

**Lemma 7.** *Let  $\mathcal{X}_C = \{ \text{vec}(\mathbf{X}) \in \mathbb{R}^{p|S|} : \text{Tr}(\mathbf{X}\mathbf{Q}\mathbf{X}^T) \leq 2C, \|\text{vec}(\mathbf{X})\|_\infty \leq B \}$ , then its covering entropy under the infinity norm is upper bounded by*

$$\log \mathcal{N}(\mathcal{X}_C, \nu) \leq p \sum_{u=1}^{|S|} \max \left( \log \left( \frac{\|\mathbf{P}\|_\infty \min(\sqrt{C\sigma_u^{-1}}, B\|\mathbf{P}\|_1)}{\sqrt{2\nu}} \right), 0 \right),$$

where  $\mathcal{N}(\mathcal{X}_C, \nu)$  is the covering number of  $\mathcal{X}_C$  under  $\|\cdot\|_\infty$ ,  $\sigma_u = \Sigma_{uu}$  is the  $u$ -th largest eigenvalue of  $\mathbf{Q}$ ,  $\mathbf{P}$  is an orthogonal matrix whose columns are the eigenvectors of  $\mathbf{Q}$ ,  $\Sigma$  is a diagonal matrix whose entries are the eigenvalues of  $\mathbf{Q}$ , and  $\mathbf{Q} = \mathbf{P}\Sigma\mathbf{P}^T$ .

**Proof of Lemma 7.** Let  $\mathbf{Z} = \mathbf{X}\mathbf{P}$ , and  $\mathbf{X} = \mathbf{Z}\mathbf{P}^T$ , and  $\mathcal{Z}_C = \{\text{vec}(\mathbf{Z}) \in \mathbb{R}^{p|S|} : \text{Tr}(\mathbf{Z}\Sigma\mathbf{Z}^T) \leq 2C, \|\text{vec}(\mathbf{Z})\|_\infty \leq B\|\mathbf{P}\|_1\}$ , and  $(B_i)_{i=1}^N$  are balls centered at  $\text{vec}(\mathbf{C}_i)$  with radius  $\nu$  to cover  $\mathcal{Z}_C$ , then we define  $\tilde{B}_i = \{(\mathbf{P} \otimes \mathbf{I}_p)\text{vec}(\mathbf{Z}) : \text{vec}(\mathbf{Z}) \in B_i\}$  with center  $\text{vec}(\tilde{\mathbf{C}}_i) = (\mathbf{P} \otimes \mathbf{I}_p)\text{vec}(\mathbf{C}_i)$ , then for any  $\text{vec}(\mathbf{X}) \in \tilde{B}_i$ , there exists  $\text{vec}(\mathbf{Z}) = (\mathbf{P} \otimes \mathbf{I}_p)^\top \text{vec}(\mathbf{X})$ , satisfying  $\text{Tr}(\mathbf{Z}\Sigma\mathbf{Z}^T) \leq 2C$ , and  $\|\text{vec}(\mathbf{Z})\|_\infty = \|(\mathbf{P} \otimes \mathbf{I}_p)^\top \text{vec}(\mathbf{X})\|_\infty \leq \|\mathbf{P}^\top\|_\infty \|\text{vec}(\mathbf{X})\|_\infty \leq B\|\mathbf{P}\|_1$ , thus there exists a ball  $\in B_i$  centered at  $\mathbf{C}_i$ , such that  $\text{vec}(\mathbf{Z}) \in B_i$  and

$$\begin{aligned} \|\text{vec}(\mathbf{X}) - \text{vec}(\tilde{\mathbf{C}})\|_\infty &\leq \|(\mathbf{P} \otimes \mathbf{I}_p)(\text{vec}(\mathbf{Z}) - \text{vec}(\mathbf{C}_i))\|_\infty \\ &= \|\mathbf{P} \otimes \mathbf{I}_p\|_\infty \|\text{vec}(\mathbf{Z}) - \text{vec}(\mathbf{C}_i)\|_\infty \leq \nu \|\mathbf{P}\|_\infty, \end{aligned}$$

where  $\otimes$  is the Kronecker product. Thus, the covering entropy of  $\mathcal{X}_C$  is upper bounded by the covering entropy of  $\mathcal{Z}_C$ , that is,  $\log \mathcal{N}(\mathcal{X}_C, \nu) \leq \log \mathcal{N}(\mathcal{Z}_C, \nu/\|\mathbf{P}\|_\infty)$ , and it suffices to compute the covering number of  $\mathcal{Z}_C$ . Note that

$$\frac{1}{2C} \text{Tr}(\mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^T) = \frac{1}{2} \sum_{u=1}^{|S|} \sum_{j=1}^p \sigma_u z_{uj}^2 / C = \sum_{u=1}^{|S|} \sum_{j=1}^p \left( \frac{z_{uj}}{\sqrt{2C/\sigma_u}} \right)^2 \leq 1, \quad (2)$$

which stands for an ellipse in  $\mathbb{R}^{p|S|}$ , with axles  $\{\sqrt{2C/\sigma_u}\}_{u=1}^{|S|}$ . Note that  $\|\text{vec}(\mathbf{Z})\|_\infty \leq B\|\mathbf{P}\|_1$ , by partitioning the axle in  $\min(\sqrt{2C/\sigma_u}, B\|\mathbf{P}\|_1)/(2\nu)$  parts, then the radius for each segment will less or equal than  $\nu$ . Thus,

$$\log \mathcal{N}(\mathcal{Z}_C, \nu) \leq p \sum_{u=1}^{|S|} \max \left( \log \left( \frac{\min(\sqrt{C\sigma_u^{-1}}, B\|\mathbf{P}\|_1)}{\sqrt{2\nu}} \right), 0 \right).$$

The desired result then follows.  $\square$

**The Armijo search in Algorithm 2.** Let  $\Theta^{(k+1)}(\gamma^{(k)})$  denote  $(\mathbf{A}_m^{(k+1)}, \mathbf{b}_m^{(k+1)})_{m=1}^M$  in



(23). Let  $\mathbf{X}^{(k+1)}(\gamma^{(k)})$  be the embedding in (22) with step size  $\gamma^{(k)}$ ,  $V_n^{(k)} = V_n(\boldsymbol{\Theta}^{(k)}, \mathbf{X}^{(k)})$  and  $\gamma^{(k)} = 2^{-t(k)}$  is determined with the Armijo search with,

$$\begin{aligned} t(k) &= \min \left\{ j \in \mathbb{Z}^+ : V_n(\boldsymbol{\Theta}^{(k+1)}(2^{-j}), \mathbf{X}^{(k+1)}(2^{-j})) \right. \\ &\quad \left. \leq V_n^{(k)} - \phi 2^{-j} \text{vec}\left(\frac{\partial V_n^{(k)}}{\partial \boldsymbol{\Theta}^{(k)}}\right)^\top \text{vec}\left(\frac{\partial V_n^{(k)}}{\partial \boldsymbol{\Theta}^{(k)}}\right) - \phi 2^{-j} \text{vec}\left(\frac{\partial V_n^{(k)}}{\partial \mathbf{X}^{(k)}}\right)^\top \text{vec}(\mathbf{X}^{(k)} - \mathbf{Z}^{(k)}) \right\}, \end{aligned}$$

for some  $\phi \in (0, 1)$ ,  $\frac{\partial V_n^{(k)}}{\partial \boldsymbol{\Theta}^{(k)}}$  and  $\frac{\partial V_n^{(k)}}{\partial \mathbf{X}^{(k)}}$  are given in (20) and (21).

## References

- [1] Jochen Gorski, Frank Pfeuffer, and Kathrin Klamroth. Biconvex sets and optimization with biconvex functions: a survey and extensions. *Mathematical Methods of Operations Research*, 66(3):373–407, 2007.
- [2] A. N. Iusem. On the convergence properties of the projected gradient method for convex optimization. *Computational and Applied Mathematics*, 22(1):37–52, 2003.
- [3] Xiaotong Shen, Lifeng Wang, et al. Generalization error for multi-class margin classification. *Electronic Journal of Statistics*, 1:307–330, 2007.
- [4] PK Suetin, Alexandra I Kostrikin, and Yu I Manin. *Linear algebra and geometry*. CRC Press, 1989.
- [5] Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, 109(3):475–494, 2001.
- [6] Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114, 2017.