

## APPENDIX

## A PRELIMINARY

**Theorem A.1** (Theorem 1.1 of (Tropp, 2012)). *Consider a finite sequence  $\{X_k\}$  of independent, random, self-adjoint matrices of dimension  $\bar{d}$ . Assume that each random matrix satisfies*

$$X_k \succeq 0 \quad \text{and} \quad \lambda(X_k) \leq R.$$

*Define  $\mu_{\min} = \lambda_{\min}(\sum_k \mathbb{E}[X_k])$  and  $\mu_{\max} = \lambda_{\max}(\sum_k \mathbb{E}[X_k])$ . Then*

$$\Pr \left\{ \lambda_{\min} \left( \sum_k X_k \right) \leq (1 - \delta) \mu_{\min} \right\} \leq \bar{d} \left( \frac{e^{-\delta}}{(1 - \delta)^{1-\delta}} \right)^{\mu_{\min}/R} \quad \text{for } \delta \in [0, 1], \text{ and} \quad (4)$$

$$\Pr \left\{ \lambda_{\max} \left( \sum_k X_k \right) \geq (1 + \delta) \mu_{\max} \right\} \leq \bar{d} \left( \frac{e^{-\delta}}{(1 + \delta)^{1+\delta}} \right)^{\mu_{\max}/R} \quad \text{for } \delta \geq 0 \quad (5)$$

**Lemma A.2** (Sherman-Morrison Formula). *Let  $A \in \mathbb{R}^{n \times n}$  be an invertible matrix, and  $u, v \in \mathbb{R}^n$ . Suppose that  $1 + v^\top A^{-1} u \neq 0$ . Then it holds that*

$$(A + uv^\top)^{-1} = A^{-1} - \frac{A^{-1}uv^\top A^{-1}}{1 + v^\top A^{-1}u}.$$

## B TWO LAYER NEURAL NETWORK

**Definition B.1** (Two layer neural network). *We define two layer neural network as follows*

$$f_{\text{nn}}(W, a, x) := \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \phi(w_r^\top x) \in \mathbb{R}$$

where  $x \in \mathbb{R}^d$  is the input,  $w_r \in \mathbb{R}^d, r \in [m]$  is the weight vector of the first layer,  $W = [w_1, \dots, w_m] \in \mathbb{R}^{d \times m}, a_r \in \mathbb{R}, r \in [m]$  is the output weight,  $a = [a_1, \dots, a_m]^\top$  and  $\phi(\cdot)$  is the non-linear activation function. Here we consider only training the first layer  $W$  with fixed  $a$ , so we also write  $f_{\text{nn}}(W, x) = f_{\text{nn}}(W, a, x)$ . Given training data matrix  $X = [x_1, \dots, x_n] \in \mathbb{R}^{n \times d}$  and labels  $Y = [y_1, \dots, y_n] \in \mathbb{R}^n$ , we denote  $f_{\text{nn}}(W, X) = [f_{\text{nn}}(W, x_1), \dots, f_{\text{nn}}(W, x_n)]^\top \in \mathbb{R}^n$ .

**Definition B.2.** *Given  $\mathcal{F}_{\text{nn}}$  and the distribution  $D$ , let  $\{v_1, \dots, v_{\bar{d}}\}$  be a fixed orthonormal basis of  $\mathcal{F}_{\text{nn}}$ , where inner products are taken under the distribution  $D$ , i.e.,*

$$\begin{aligned} \mathbb{E}_{x \sim D} [v_i(x) \cdot v_j(x)] &= 1, \forall i = j \in [\bar{d}] \\ \left| \mathbb{E}_{x \sim D} [v_i(x) \cdot v_j(x)] \right| &\leq \rho, \forall i \neq j \in [\bar{d}] \end{aligned}$$

Furthermore, for any function  $h \in \mathcal{F}_{\text{nn}}$ , there exists  $\alpha(h) := (\alpha(h)_1, \dots, \alpha(h)_{\bar{d}})$  under the basis  $(v_1, \dots, v_{\bar{d}})$  such that

$$h(x) = \sum_{i=1}^{\bar{d}} \alpha(h)_i \cdot v_i(x)$$

**Remark.** *Note that in for linear function family  $\mathcal{F}_{\text{nn}}$ , we know that  $\bar{d} = d$ . However, for the neural network function family  $\mathcal{F}_{\text{nn}}$ , we will have  $\bar{d} \gg d$ .*

**Definition B.3.** *Given distribution  $D$  and  $h \in \mathcal{F}_{\text{nn}}$ , we define  $\|h(x)\|_D$  as*

$$\|h(x)\|_D^2 := \mathbb{E}_{x \sim D} [|h(x)|^2].$$

**Claim B.4.** *If the activation  $\phi$  in a neural network  $f_{\text{nn}}$  satisfies the following conditions,*

- $\phi^{(10d + \log(1/\varepsilon_0)/\log(d))}(x)$  exists and is continuous.
- $\phi^{(10d + \log(1/\varepsilon_0)/\log(d))}(x) \leq 1, x \in \mathbb{R}$ .

then there exists  $\{v_1, \dots, v_{\bar{d}}\}$  forms a fixed  $\rho$ -nearly orthonormal basis of  $\mathcal{F}_{\text{nn}}$  where

$$\bar{d} \leq \binom{10d + \log(1/\varepsilon_0)/\log(d)}{d}.$$

Furthermore, for any  $W \in \mathbb{R}^{d \times m}$ , there exists  $h \in \mathcal{F}_{\text{nn}}$  such that

$$\|h(x) - f_{\text{nn}}(W, x)\|_D^2 \leq \varepsilon.$$

Besides, for any  $h \in \mathcal{F}_{\text{nn}}$ , there exists  $W \in \mathbb{R}^{d \times m}$  such that

$$\|h(x) - f_{\text{nn}}(W, x)\|_D^2 \leq \varepsilon.$$

*Proof.* For any activation  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  and input  $z_r = w_r^\top x \in \mathbb{R}$ ,  $f$  can be expand by Taylor's theorem

$$\phi(z_r) = \phi(0) + \phi'(0)z_r + \frac{\phi''(0)}{2!}z_r^2 + \dots + \frac{\phi^{(k)}(0)}{k!}z_r^k + \frac{\phi^{(k+1)}(\xi)}{(k+1)!}z_r^{k+1}$$

where  $\xi \in [0, z_r]$ .

It's natural to consider the  $W$  bound, eg NTK regime. In NTK regime, we have with a high probability that

$$z_r = w_r^\top x \leq \|x\|_1 \leq \sqrt{d}\|x\|_2 \leq \sqrt{d}.$$

We can claim that

$$\frac{\phi^{(k+1)}(\xi)}{(k+1)!}z_r^{k+1} \leq \left(\frac{e}{k+1}z_r\right)^{k+1} \leq \varepsilon_0.$$

where the first step follows from  $\phi^{(k+1)}(\xi) \leq 1$ , the second step follows from  $k \geq (e\sqrt{d})^2 + \log(1/\varepsilon_0)/\log(e\sqrt{d})$ .

Besides, taking the first  $k+1$  terms in the Taylor's theorem. Our neural network  $f_{\text{nn}}$  can be seen as a polynomial with  $d$  variable and at most  $k$  degree. So, our neural network  $f_{\text{nn}}$  can be seen as a polynomial with at most  $\binom{k+d}{d}$  terms. As any polynomial can be orthonormal decompose, we have that

$$\bar{d} \leq \binom{k+d}{d} \leq \binom{10d + \log(1/\varepsilon_0)/\log(d)}{d}.$$

where the second step follows from  $d \geq e$ . □

**Claim B.5.** We will claim that polynomial, ReLU, Sigmoid, and Swish hold the condition mentioned in Claim B.4.

**Lemma B.6.** Let  $A \in \mathbb{R}^{\bar{d} \times \bar{d}}$  be defined as

$$\begin{aligned} A(i, j) &= 1, \forall i = j \in [\bar{d}] \\ |A(i, j)| &\leq \rho, \forall i \neq j \in [\bar{d}]. \end{aligned}$$

We will claim that

$$\lambda_{\max}(A) \leq 1 + \bar{d}\rho, \lambda_{\min}(A) \geq 1 - \bar{d}\rho.$$

*Proof.* First, we can lower bound  $\lambda_{\min}(A)$  as follows

$$\begin{aligned} \lambda_{\min}(A) &\geq \lambda_{\min}(I) - \|I - A\|_2 \\ &\geq 1 - \|I - A\|_F \\ &\geq 1 - \bar{d}\rho \end{aligned}$$

Second, we can upper bound  $\lambda_{\max}(A)$  as follows

$$\begin{aligned} \lambda_{\max}(A) &= \|A\|_2 \\ &\leq \|I\|_2 + \|A - I\|_2 \\ &\leq 1 + \|A - I\|_F \\ &\leq 1 + \bar{d}\rho \end{aligned}$$

□

## C NOTATIONS

**Claim C.1.** For any function  $h \in \mathcal{F}_{\text{nn}}$ , we have

$$C_l \|\alpha(h)\|_2^2 \leq \|h\|_D^2 \leq C_r \|\alpha(h)\|_2^2$$

where  $C_l := 1 - \rho$  and  $C_r := 1 + \rho(\bar{d} - 1)$ .

*Proof.* We can rewrite  $\|h\|_D^2$  as follows:

$$\begin{aligned} \mathbb{E}_{x \sim D} \left[ \left| \sum_{i=1}^{\bar{d}} \alpha(h)_i \cdot v_i(x) \right|^2 \right] &= \sum_{i=1}^{\bar{d}} \mathbb{E}_{x \sim D} [|\alpha(h)_i \cdot v_i(x)|^2] + 2 \sum_{1 \leq i < j \leq \bar{d}} \mathbb{E}_{x \sim D} [\alpha(h)_i \alpha(h)_j \cdot v_i(x) v_j(x)] \\ &= \sum_{i=1}^{\bar{d}} |\alpha(h)_i|^2 + 2\rho \cdot \sum_{1 \leq i < j \leq \bar{d}} \alpha(h)_i \alpha(h)_j \\ &= (1 - \rho) \|\alpha(h)\|_2^2 + \rho \left( \sum_{i=1}^{\bar{d}} \alpha(h)_i \right)^2 \end{aligned}$$

We can provide an upper bound,

$$\begin{aligned} \mathbb{E}_{x \sim D} \left[ \left| \sum_{i=1}^{\bar{d}} \alpha(h)_i \cdot v_i(x) \right|^2 \right] &= (1 - \rho) \|\alpha(h)\|_2^2 + \rho \left( \sum_{i=1}^{\bar{d}} \alpha(h)_i \right)^2 \\ &\leq (1 - \rho) \|\alpha(h)\|_2^2 + \rho \bar{d} \|\alpha(h)\|_2^2 \\ &= (1 - \rho + \rho \bar{d}) \|\alpha(h)\|_2^2 \\ &= C_r \|\alpha(h)\|_2^2 \end{aligned}$$

We can provide a lower bound

$$\begin{aligned} \mathbb{E}_{x \sim D} \left[ \left| \sum_{i=1}^{\bar{d}} \alpha(h)_i \cdot v_i(x) \right|^2 \right] &= (1 - \rho) \|\alpha(h)\|_2^2 + \rho \left( \sum_{i=1}^{\bar{d}} \alpha(h)_i \right)^2 \\ &\geq (1 - \rho) \|\alpha(h)\|_2^2 \\ &= C_l \|\alpha(h)\|_2^2 \end{aligned}$$

Thus, we complete the proof.  $\square$

Now, we define condition number. Previous work only consider linear cases and we generalize it to NN.

**Definition C.2.** For any distribution  $D'$  over the domain  $G$  and any function  $h : G \rightarrow \mathbb{R}$ , let  $h^{(D')}(x) = \sqrt{\frac{D(x)}{D'(x)}} \cdot h(x)$  such that  $\mathbb{E}_{x \sim D'} [ |h^{(D')}(x)|^2 ] = \mathbb{E}_{x \sim D'} \left[ \frac{D(x)}{D'(x)} |h(x)|^2 \right] = \mathbb{E}_{x \sim D} [ |h(x)|^2 ]$ . When the neural network function  $\mathcal{F}_{\text{nn}}$  and  $D$  is clear, we use  $K_{D'}$  to denote the condition number of sampling from  $D'$ , i.e.,

$$K_{D'} = \sup_x \left\{ \sup_{h \in \mathcal{F}_{\text{nn}}} \left\{ \frac{|h^{(D')}(x)|^2}{\|h^{(D')}\|_{D'}^2} \right\} \right\} = \sup_x \left\{ \frac{D(x)}{D'(x)} \cdot \sup_{h \in \mathcal{F}_{\text{nn}}} \left\{ \frac{|h(x)|^2}{\|h\|_D^2} \right\} \right\}.$$

**Definition C.3.** For any distribution  $D'$  over the domain  $G$  and any function  $h : G \rightarrow \mathbb{R}$ . When the neural network function  $\mathcal{F}_{\text{nn}}$  and  $D$  is clear, we use  $K_{\alpha, D'}$  to denote the  $\alpha$ -condition number of sampling from  $D'$ , i.e.,

$$K_{\alpha, D'} = \sup_x \left\{ \frac{D(x)}{D'(x)} \cdot \sup_{h \in \mathcal{F}_{\text{nn}}} \left\{ \frac{|h(x)|^2}{\|\alpha(h)\|_2^2} \right\} \right\}.$$

**Definition C.4.** Given  $\mathcal{F}_{\text{nn}}$  and underlying distribution  $D$ , let  $P$  be a random sampling procedure that terminates in  $k$  iterations ( $k$  is not necessarily fixed) and provides a coefficient  $\alpha_i$  and a distribution  $D_i$  to sample  $x_i \sim D_i$  in every iteration  $i \in [k]$ .

We say  $P$  is an  $\varepsilon$ -importance sampling procedure if it satisfies the following two properties:

1. Let  $v_1, \dots, v_{\bar{d}}$  of  $\mathcal{F}_{nn}$  under  $D$  be defined as in Definition B.2. With probability 0.9, the matrix  $A(i, j) = \sqrt{u_i} \cdot v_j(x_i) \in \mathbb{R}^{k \times \bar{d}}$  has  $\lambda(A^* A) \in [\frac{3}{4}, \frac{5}{4}]$ .
2. The coefficients always have  $\sum_{i=1}^k \beta_i \leq \frac{5}{4}$  and  $\beta_i \cdot K_{\alpha, D_i} \leq \epsilon/2$ .

**Claim C.5.** We will claim that

$$C_l K_{D'} \leq K_{\alpha, D'} \leq C_r K_{D'}$$

*Proof.* Since, we have that

$$C_l \frac{1}{\|h\|_D^2} \leq \frac{1}{\|\alpha(h)\|_2^2} \leq C_r \left\{ \frac{1}{\|h\|_D^2} \right\},$$

we can claim that

$$C_l \cdot \sup_x \left\{ \frac{D(x)}{D'(x)} \cdot \sup_{h \in \mathcal{F}_{nn}} \left\{ \frac{|h(x)|^2}{\|h\|_D^2} \right\} \right\} \leq \sup_x \left\{ \frac{D(x)}{D'(x)} \cdot \sup_{h \in \mathcal{F}_{nn}} \left\{ \frac{|h(x)|^2}{\|\alpha(h)\|_2^2} \right\} \right\} \leq C_r \cdot \sup_x \left\{ \frac{D(x)}{D'(x)} \cdot \sup_{h \in \mathcal{F}_{nn}} \left\{ \frac{|h(x)|^2}{\|h\|_D^2} \right\} \right\}$$

□

**Definition C.6.** Given a importance sampling procedure  $P$ , we say the output of  $P$  is good only if the samples  $x_i$  with weights  $u_i = \beta_i \cdot D(x_i)/D_i(x_i)$  satisfy the first property in Definition C.4. Given a joint distribution  $(D, Y)$  and an execution of a importance sampling procedure  $P$  with  $x_i \sim D_i$  and  $u_i$  of each  $i \in [k]$ , let the  $\tilde{f}$  be defined as

$$\tilde{f} = \arg \min_{h \in \mathcal{F}_{nn}} \left\{ \sum_{i=1}^k u_i \cdot |h(x_i) - y_i|^2 \right\}$$

by querying  $y_i \sim (Y|x_i)$  for each point  $x_i$ .

## D RECOVERY GUARANTEE FOR IMPORTANCE SAMPLES

### D.1 PRELIMINARY

We state a tool from prior work,

**Lemma D.1** (Lemma 4.3 of (Chen & Price, 2019)). Let  $P$  be a random sampling procedure terminating in  $k$  iterations ( $k$  is not necessarily fixed) that in every iteration  $i$ , it provides a coefficient  $\beta_i$  and a distribution  $D_i$  to sample  $x_i \sim D_i$ . Let the weight  $u_i = \beta_i \cdot \frac{D(x_i)}{D_i(x_i)}$  and  $A \in \mathbb{R}^{k \times \bar{d}}$  denote the matrix  $A(i, j) = \sqrt{u_i} \cdot v_j(x_i)$ . Then for  $f = \arg \min_{h \in \mathcal{F}_{nn}} \mathbb{E}_{(x, y) \sim (D, Y)} [|y - h(x)|^2]$ ,

$$\mathbb{E}_P \left[ \|A^*(\vec{y}_u - \vec{f}_{S,u})\|_2^2 \right] \leq \sup_P \left\{ \sum_{i=1}^k \beta_i \right\} \cdot \max_j \left\{ \beta_j \cdot K_{\alpha, D_j} \right\}_{(x, y) \sim (D, Y)} \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2],$$

where  $K_{\alpha, D_i}$  is the condition number for samples from  $D_i$ :  $K_{\alpha, D_i} = \sup_x \left\{ \frac{D(x)}{D_i(x)} \cdot \sup_{v \in \mathcal{F}} \left\{ \frac{|v(x)|^2}{\|\alpha(v)\|_2^2} \right\} \right\}$ .

### D.2 ANALYSIS FOR IMPORTANCE SAMPLING PROCEDURE

**Lemma D.2.** For any  $\epsilon \in (0, 1)$ , given  $S = (x_1, \dots, x_k) \subset \mathbb{R}^d$  and their weights  $(u_1, \dots, u_k) \subset \mathbb{R}_{\geq 0}$ , let  $A$  be the  $k \times \bar{d}$  matrix defined as

$$A(i, j) := \sqrt{u_i} \cdot v_j(x_i).$$

Then,

$$\lambda(A^* A) \in [1 - \epsilon, 1 + \epsilon].$$

can imply that

$$\|h\|_{S,u}^2 := \sum_{j=1}^k u_j \cdot |h(x_j)|^2 \in \left[ \frac{1-\epsilon}{C_r}, \frac{1+\epsilon}{C_l} \right] \cdot \|h\|_D^2 \quad \text{for every } h \in \mathcal{F}_{nn}$$

*Proof.* Notice that

$$A \cdot \alpha(h) = (\sqrt{u_1} \cdot h(x_1), \dots, \sqrt{u_k} \cdot h(x_k)). \quad (6)$$

We can show

$$\begin{aligned} \|h\|_{S,u}^2 &= \sum_{i=1}^k u_i \cdot |h(x_i)|^2 \\ &= \|A \cdot \alpha(h)\|_2^2 \\ &= \alpha(h)^* \cdot (A^* \cdot A) \cdot \alpha(h) \\ &\in [\lambda_{\min}(A^* \cdot A), \lambda_{\max}(A^* \cdot A)] \cdot \|\alpha(h)\|_2^2 \\ &\subseteq [1 - \varepsilon, 1 + \varepsilon] \cdot \|\alpha(h)\|_2^2 \\ &\subseteq [\frac{1 - \varepsilon}{C_r}, \frac{1 + \varepsilon}{C_l}] \cdot \|h\|_D^2 \end{aligned}$$

where the first step follows from the definition of the norm, the second step follows from the Eq. (6), the last step follows from  $C_l \|\alpha(h)\|_2^2 \leq \|h\|_D^2 \leq C_r \|\alpha(h)\|_2^2$ .  $\square$

For any  $\varepsilon > 0$ , given  $S = (x_1, \dots, x_k) \subset \mathbb{R}^d$  and their weights  $(u_1, \dots, u_k) \subset \mathbb{R}_{\geq 0}$ . Let  $A$  be the  $k \times \bar{d}$  matrix defined as

$$A(i, j) := \sqrt{u_i} \cdot v_j(x_i).$$

For any  $h \in \mathcal{F}_{nn}$ , let  $\|h\|_{S,u}^2$  be defined as

$$\|h\|_{S,u}^2 := \sum_{j=1}^k u_j \cdot |h(x_j)|^2$$

Then, for any  $h \in \mathcal{F}_{nn}$

$$A \cdot \alpha(h) = (\sqrt{u_1} \cdot h(x_1), \dots, \sqrt{u_k} \cdot h(x_k)). \quad (7)$$

We consider the calculation of the  $\tilde{f}$ . Given the weights  $(u_1, \dots, u_k)$  on  $(x_1, \dots, x_k)$  and labels  $(y_1, \dots, y_k)$ , let  $\vec{y}_u$  denote the vector of weighted labels  $(\sqrt{u_1} \cdot y_1, \dots, \sqrt{u_k} \cdot y_k)$ . From Eq. (6), the empirical distance satisfied that for any  $h \in \mathcal{F}_{nn}$

$$\begin{aligned} \|h(x_i) - y_i\|_{S,u}^2 &= \sum_{i=1}^k u_i |h(x_i) - y_i|^2 \\ &= \|A \cdot \alpha(h) - \vec{y}_u\|_2^2 \end{aligned}$$

where the second step follows from Eq. (6) and the definition of  $\vec{y}_u$ .

Let

$$\begin{aligned} \tilde{f} &= \arg \min_{h \in \mathcal{F}_{nn}} \{\|h(x_i) - y_i\|_{S,u}\} \\ &= \arg \min_{h \in \mathcal{F}_{nn}} \{\|A \cdot \alpha(h) - \vec{y}_u\|_2\} \end{aligned}$$

Then,

$$\alpha(\tilde{f}) = (A^* \cdot A)^{-1} \cdot A^* \cdot \vec{y}_u \text{ and } \tilde{f} = \sum_{i=1}^{\bar{d}} \alpha(\tilde{f})_i \cdot v_i.$$

Let

$$f = \arg \min_{h \in \mathcal{F}_{nn}} \{ \mathbb{E}_{(x,y) \sim (D,Y)} [|h(x) - y|^2] \}$$

Finally, we consider the distance between  $f$  and  $\tilde{f}$ . For convenience, let  $\vec{f}_u = (\sqrt{u_1} \cdot f(x_1), \dots, \sqrt{u_k} \cdot f(x_k))$ . Because  $f \in \mathcal{F}_{\text{nn}}$  and Eq. (6), we can claim that,

$$\alpha(f) = (A^* \cdot A)^{-1} \cdot A^* \cdot \vec{f}_u.$$

This implies

$$\|\tilde{f} - f\|_D^2 \in [C_l, C_r] \cdot \|\alpha(\tilde{f}) - \alpha(f)\|_2^2 = [C_l, C_r] \cdot \|(A^* \cdot A)^{-1} \cdot A^* \cdot (\vec{y}_u - \vec{f}_u)\|_2^2,$$

where the first step follows from  $\|h\|_D^2 \in [C_l, C_r] \|\alpha(h)\|_2^2$  and  $\alpha(f - g) = \alpha(f) - \alpha(g)$ .

**Theorem D.3.** *Given a neural network function family  $\mathcal{F}_{\text{nn}}$ , joint distribution  $(D, Y)$ , and  $\varepsilon > 0$ , let  $P$  be an  $\varepsilon$ -importance sampling procedure for  $\mathcal{F}_{\text{nn}}$  and  $D$ , and let  $f = \arg \min_{h \in \mathcal{F}} \mathbb{E}_{(x,y) \sim (D,Y)} [|y - h(x)|^2]$ . Then the  $\tilde{f}$  of a good output of  $P$  satisfies*

$$\|f - \tilde{f}\|_D^2 \leq \epsilon \cdot \mathbb{E}_{(x,y) \sim (D,Y)} [|y - f(x)|^2] \text{ in expectation.}$$

*Proof.* We assume the first property  $\lambda(A^* \cdot A) \in [1 - 1/4, 1 + 1/4]$  from Definition C.6. On the other hand,

$$\mathbb{E}_P [\|A^* \cdot (\vec{y}_u - \vec{f}_u)\|_2^2] \leq \epsilon \cdot \mathbb{E}_{(x,y) \sim (D,Y)} [|y - f(x)|^2]$$

from Lemma D.1 where  $P$  is an random sampling procedure. Conditioned on the first property, we know

$$\begin{aligned} \mathbb{E}_{P'} [\|A^* \cdot (\vec{y}_u - \vec{f}_u)\|_2^2] &\leq \frac{1}{0.9} \mathbb{E}_P [\|A^* \cdot (\vec{y}_u - \vec{f}_u)\|_2^2] \\ &\leq \frac{\varepsilon}{0.9} \cdot \mathbb{E}_{(x,y) \sim (D,Y)} [|y - f(x)|^2] \end{aligned}$$

where  $P'$  is the event where importance sampling procedure  $P$  good executed.

This implies

$$\begin{aligned} \mathbb{E}_{P'} [\|(A^* \cdot A)^{-1} \cdot A^* \cdot (\vec{y}_u - \vec{f}_u)\|_2^2] &\leq \mathbb{E}_{P'} [\lambda_{\min}(A^* \cdot A)^{-1} \cdot \|A^* \cdot (\vec{y}_u - \vec{f}_u)\|_2^2] \\ &\leq 2\epsilon \cdot \mathbb{E}_{(x,y) \sim (D,Y)} [|y - f(x)|^2] \end{aligned}$$

where the second step follows from the first property in the definition of a importance sampling procedure  $P$  and  $P'$  is a good output of  $P$ .  $\square$

## E A LINEAR-SAMPLE ALGORITHM FOR KNOWN DISTRIBUTION

### E.1 PRELIMINARY

We state several tools from prior work.

**Lemma E.1** (Lemma 3.3 in (Batson et al., 2012)). *For any  $j \in [k]$ ,  $\lambda(B_j) \in (l_j, r_j)$ .*

**Lemma E.2** (A combination of Claim 5.5 and Lemma 5.6 in (Chen & Price, 2019)). *After exiting the while loop in Procedure RANDOMIZEDBSS, we always have*

1.  $r_k - l_k \leq 9\bar{d}/\gamma$ .
2.  $(1 - \frac{0.5\gamma^2}{\bar{d}}) \cdot \sum_{j=1}^k \frac{\gamma}{\phi_j} \leq \text{mid} \leq \sum_{j=1}^k \frac{\gamma}{\phi_j}$ .
3. If  $\frac{r_k}{l_k} \leq 1 + 8\gamma$ , then  $\lambda(A^* \cdot A) \in (1 - 5\gamma, 1 + 5\gamma)$ .

**Lemma E.3** (Lemma 5.1 in (Chen & Price, 2019)). *Given any dimension  $\bar{d}$  linear space  $\mathcal{F}_{\text{nn}}$ , any distribution  $D$  over the domain of  $\mathcal{F}_{\text{nn}}$ , and any  $\varepsilon > 0$ , there exists an  $\varepsilon$ -importance sampling procedure that terminates in  $O(\bar{d}/\varepsilon)$  rounds with probability 0.99.*

**Algorithm 2** A importance sampling procedure based on Randomized Sampling

---

```

1: procedure RANDOMIZEDSAMPLING( $\mathcal{F}_{nn}, D, \epsilon$ )
2:   Find an  $\rho$ -nearly orthonormal basis  $v_1, \dots, v_{\bar{d}}$  of  $\mathcal{F}_{nn}$  under  $D$ ;
3:    $\gamma \leftarrow \sqrt{\epsilon}/C_0$ ;
4:    $\text{mid} \leftarrow (4\bar{d}/\gamma)/(1/(1-\gamma) - 1/(1+\gamma))$ ;
5:    $B_0 \leftarrow 0$ ;
6:    $l_0 \leftarrow -2\bar{d}/\gamma$ ;
7:    $r_0 \leftarrow 2\bar{d}/\gamma$ ;
8:    $j \leftarrow 0$ ;
9:   while  $r_{j+1} - l_{j+1} < 8\bar{d}/\gamma$  do;
10:     $\Phi_j \leftarrow \text{tr}[(r_j I - B_j)^{-1}] + \text{tr}[(B_j - l_j I)^{-1}]$ ;
11:     $D_j(x) \leftarrow D(x) \cdot \left( v(x)^\top (r_j I - B_j)^{-1} v(x) + v(x)^\top (B_j - l_j I)^{-1} v(x) \right) / \Phi_j$ ;
12:    Sample  $x_j \sim D_j$ ;
13:     $s_j \leftarrow \gamma \cdot D(x) / (\Phi_j \cdot D_j(x))$ ;
14:     $B_{j+1} \leftarrow B_j + s_j \cdot v(x_j) v(x_j)^\top$ ;
15:     $r_{j+1} \leftarrow r_j + \gamma / (\Phi_j(1-\gamma))$ ;
16:     $l_{j+1} \leftarrow l_j + \gamma / (\Phi_j(1+\gamma))$ ;
17:     $j \leftarrow j + 1$ ;
18:  end while
19:   $k \leftarrow j$ ;
20:  for  $j \in [k]$  do
21:     $\beta_j \leftarrow \gamma / (\Phi_j \cdot \text{mid})$ ,
22:     $u_j \leftarrow s_j / \text{mid}$ ,
23:  end for
24:  Output  $x, D, u, \beta$ 
25: end procedure

```

---

## E.2 ANALYSIS FOR OUR RANDOMIZED SAMPLING

Our results in Lemma E.4, Lemma E.5, Lemma E.6, and E.7 are different than results in (Allen-Zhu et al., 2015; Lee & Sun, 2018). First, they consider a potential function  $\text{tr}[(r_j I - B_j)^{-q}] + \text{tr}[(B_j - l_j I)^{-q}]$  with  $q \geq 10$ . But, in our cases,  $q = 1$ . Second, their results is for orthonormal basis but not for  $\rho$ -nearly orthonormal basis.

**Lemma E.4.** Let  $\epsilon \in (0, 1/2)$ . Suppose that  $w^\top (uI - A)^{-1} w \leq \epsilon$  and  $w^\top (A - lI)^{-1} w \leq \epsilon$ . It holds that

$$\begin{aligned} \text{tr}[(A - lI + ww^\top)^{-1}] &\leq \text{tr}[(A - lI)^{-1}] - (1 - \epsilon)w^\top (A - lI)^{-2} w \\ \text{tr}[(uI - A - ww^\top)^{-1}] &\leq \text{tr}[(uI - A)^{-1}] + (1 + 2\epsilon)w^\top (uI - A)^{-2} w. \end{aligned}$$

*Proof.* Let  $Y = A - lI$ . By the Sherman-Morrison Formula (Lemma A.2), it holds that

$$\text{tr}[(Y + ww^\top)^{-1}] = \text{tr}[Y^{-1} - \frac{Y^{-1}ww^\top Y^{-1}}{1 + w^\top Y^{-1}w}].$$

By the assumption of  $w^\top Y^{-1}w \leq \epsilon$ , we have that

$$\begin{aligned} \text{tr}[(Y + ww^\top)^{-1}] &\leq \text{tr}[Y^{-1} - \frac{Y^{-1}ww^\top Y^{-1}}{1 + \epsilon}] \\ &= \text{tr}[Y^{-1}(I - \frac{Y^{-1/2}ww^\top Y^{-1/2}}{1 + \epsilon})] \\ &\leq \text{tr}[Y^{-1}(I - (1 - \epsilon)Y^{-1/2}ww^\top Y^{-1/2})] \\ &\leq \text{tr}[Y^{-1}] - (1 - \epsilon)w^\top Y^{-2}w. \end{aligned}$$

Let  $Z = uI - A$ . By the Sherman-Morrison Formula (Lemma A.2), it holds that

$$\text{tr}[(Z - ww^\top)^{-1}] = \text{tr}[Z^{-1} + \frac{Z^{-1}ww^\top Z^{-1}}{1 - w^\top Z^{-1}w}].$$

By the assumption of  $w^\top Z^{-1}w \leq \varepsilon$ , we have that

$$\begin{aligned} \text{tr}[(Z - ww^\top)^{-1}] &\leq \text{tr}[Z^{-1} + \frac{Z^{-1}ww^\top Z^{-1}}{1 - \varepsilon}] \\ &= \text{tr}[Z^{-1}(I + \frac{Z^{-1/2}ww^\top Z^{-1/2}}{1 - \varepsilon})] \\ &\leq \text{tr}[Z^{-1}(I + (1 + 2\varepsilon)Z^{-1/2}ww^\top Z^{-1/2})] \\ &\leq \text{tr}[Z^{-1}] + (1 + 2\varepsilon)w^\top Z^{-2}w. \end{aligned}$$

where the third step follows from  $\varepsilon \leq 0.5$ . □

**Lemma E.5.** *Let  $w_i$  be defined as*

$$w_j := \sqrt{\frac{\gamma}{v(x_j)^\top (r_j I - B_j)^{-1} v(x_j) + v(x_j)^\top (B_j - l_j I)^{-1} v(x_j)}} \cdot v(x_j).$$

*Let  $\gamma \leq \min\{1/(c \log(\bar{d}/\gamma')), 1/(c(1 + \bar{d}\rho)e^2)\}$ . Then, it holds that*

$$\begin{aligned} \Pr[0 \preceq w_j w_j^\top \preceq \frac{1}{c} \cdot (r_j I - B_j)] &\geq 1 - \gamma' \\ \Pr[0 \preceq w_j w_j^\top \preceq \frac{1}{c} \cdot (B_j - l_j I)] &\geq 1 - \gamma' \end{aligned}$$

*Proof.* Let  $R_j = v(x_j)^\top (r_j I - B_j)^{-1} v(x_j) + v(x_j)^\top (B_j - l_j I)^{-1} v(x_j)$ .

We can claim that

$$\mathbb{E}_{x \sim D_j} [w_j w_j^\top] = \frac{\gamma}{\Phi_j} \mathbb{E}_{x \sim D} [v(x_j) v(x_j)^\top] \preceq \frac{\gamma(1 + \bar{d}\rho)}{\Phi_j} \cdot I$$

Let

$$z_j = (r_j I - B_j)^{-1/2} w_j.$$

It holds that

$$\begin{aligned} \text{tr}[z_j z_j^\top] &= \text{tr}[(r_j I - B_j)^{-1/2} w_j w_j^\top (r_j I - B_j)^{-1/2}] \\ &= \frac{\gamma}{R_j} \cdot \text{tr}[(r_j I - B_j)^{-1/2} v_j v_j^\top (r_j I - B_j)^{-1/2}] \\ &= \frac{\gamma}{R_j} \cdot \text{tr}[v_j^\top (r_j I - B_j)^{-1} v_j] \\ &\leq \gamma, \end{aligned}$$

and  $\lambda_{\max}(z_j z_j^\top) \leq \gamma$ . Moreover, it holds that

$$\begin{aligned} \mathbb{E}[z_j z_j^\top] &= \frac{\gamma(1 + \bar{d}\rho)}{\Phi_j} \cdot (r_j I - B_j)^{-1} \\ &\preceq \frac{\gamma}{\Phi_j} \cdot \lambda_{\max}\left(\frac{1}{r_j I - B_j}\right) \cdot I. \end{aligned}$$

This implies that

$$\lambda_{\max}(\mathbb{E}[z_j z_j^\top]) \leq \frac{\gamma(1 + \bar{d}\rho)}{\Phi_j} \cdot \lambda_{\max}\left(\frac{1}{r_j I - B_j}\right) =: \mu$$

It holds by the Matrix Chernoff Bound (Lemma A.1) that

$$\Pr[\lambda_{\max}(\mathbb{E}[z_j z_j^\top]) \geq (1 + \delta)\mu] \leq \bar{d} \cdot \left(\frac{\exp(\delta)}{(1 + \delta)^{(1 + \delta)}}\right)^{\mu/\gamma}.$$



Set  $1 + \delta$  to be

$$\begin{aligned} 1 + \delta &= 1/(c\mu) \\ &= \frac{\Phi_j}{c\gamma(1 + \bar{d}\rho)} \lambda_{\min}(r_j I - B_j) \\ &\geq \frac{1}{c\gamma(1 + \bar{d}\rho)} \end{aligned}$$

With probability at least

$$\begin{aligned} 1 - \bar{d} \cdot \left( \frac{\exp(\delta)}{(1 + \delta)^{(1 + \delta)}} \right)^{\mu/\gamma} &\geq 1 - \bar{d} \cdot \left( \frac{e}{1 + \delta} \right)^{\mu(1 + \delta)/\gamma} \\ &= 1 - \bar{d} \cdot \left( \frac{e}{1 + \delta} \right)^{1/c\gamma} \\ &\geq 1 - \bar{d} \cdot (ce\gamma(1 + \bar{d}\rho))^{1/(c\gamma)} \\ &\geq 1 - \bar{d} \cdot \exp(-1/(c\gamma)) \\ &\geq 1 - \gamma' \end{aligned}$$

where the fourth step follows from  $\gamma \leq 1/(c(1 + \bar{d}\rho)e^2)$ , the fifth step follows from  $\gamma \leq 1/(c \log(\bar{d}/\gamma'))$ .

As a result, we can prove that

$$\Pr[0 \preceq w_j w_j^\top \preceq \frac{1}{c} \cdot (r_j I - B_j)] \geq 1 - \gamma'.$$

Similarly, we can prove that

$$\Pr[0 \preceq w_j w_j^\top \preceq \frac{1}{c} \cdot (B_j - l_j I)] \geq 1 - \gamma'.$$

□

**Lemma E.6.** *It holds that*

$$\mathbb{E}_{x_j \in D_j} [\Phi_{j+1}] \leq \Phi_j.$$

*Proof.* We claim that

$$w_j w_j^\top \preceq \frac{1}{c} \cdot (r_j I - B_j) \preceq \frac{1}{c} \cdot (r_{j+1} I - B_j)$$

We apply Lemma E.4 with  $\varepsilon = 1/c$  and get that

$$\begin{aligned} &\mathbb{E}_{x_j \in D_j} [\text{tr}[(r_{j+1} I - B_{j+1})^{-1}]] \\ &\leq \text{tr}[(r_{j+1} I - B_j)^{-1}] + (1 + 2/c) \text{tr}[(r_{j+1} I - B_j)^{-2}] \mathbb{E}[w_j w_j^\top] \\ &= \text{tr}[(r_{j+1} I - B_j)^{-1}] + \frac{(1 + 2/c)\gamma}{\Phi_j} \text{tr}[(r_{j+1} I - B_j)^{-2}] \end{aligned}$$

Note that  $r_{j+1} - r_j = \gamma/(\Phi_j(1 - \gamma))$ , we define a function  $f$  by

$$f(t) = \text{tr}[(r_j + t \cdot (r_{j+1} - r_j))I - B_j]^{-1}.$$

Notice that

$$\frac{df(t)}{dt} = -\frac{\gamma}{\Phi_j(1 - \gamma)} \text{tr}[(r_j + t \cdot (r_{j+1} - r_j))I - B_j]^{-2}$$

Since  $f$  is convex, we have that

$$\left. \frac{df(t)}{dt} \right|_{t=1} \geq f(1) - f(0) = \text{tr}[(r_{j+1}I - B_j)^{-1}] - \text{tr}[(r_jI - B_j)^{-1}].$$

We can conclude that

$$\begin{aligned} \mathbb{E}_{x_j \in D_j} [\text{tr}[(r_{j+1}I - B_{j+1})^{-1}]] &\leq \text{tr}[(r_{j+1}I - B_j)^{-1}] + \frac{(1 + 2/c)\gamma}{\Phi_j} \text{tr}[(r_{j+1}I - B_j)^{-2}] \\ &\leq \text{tr}[(r_{j+1}I - B_j)^{-1}] \\ &\quad + (1 + 2/c)(1 - \gamma)(\text{tr}[(r_jI - B_j)^{-1}] - \text{tr}[(r_{j+1}I - B_j)^{-1}]) \\ &\leq \text{tr}[(r_jI - B_j)^{-1}] \end{aligned}$$

where the last step follows from  $(1 + 2/c)(1 - \gamma) \leq 1$ .

On the other hand, we can proof that

$$w_j w_j^T \preceq \frac{1}{c} \cdot (B_j - l_j I) \preceq \frac{1}{c} \cdot (B_j - l_{j+1} I)$$

We apply Lemma E.4 with  $\varepsilon = 1/c$  and get that

$$\begin{aligned} \mathbb{E}_{x_j \in D_j} [\text{tr}[(B_{j+1} - l_{j+1} I)^{-1}]] &\leq \text{tr}[(B_j - l_{j+1} I)^{-1}] - (1 - 1/c) \text{tr}[(B_j - l_{j+1} I)^{-2} \mathbb{E}[w_j w_j^T]] \\ &= \text{tr}[(B_j - l_{j+1} I)^{-1}] - \frac{(1 - 1/c)\gamma}{\Phi_j} \text{tr}[(B_j - l_{j+1} I)^{-2}] \end{aligned}$$

Note that  $l_{j+1} - l_j = \gamma/(\Phi_j(1 + \gamma))$ , we define a function  $g$  by

$$g(t) = \text{tr}[(B_j - (l_j + t \cdot (l_{j+1} - l_j))I)^{-1}].$$

Notice that

$$\left. \frac{dg(t)}{dt} \right|_{t=1} = \frac{\gamma}{\Phi_j(1 + \gamma)} \text{tr}[(B_j - (l_j + t \cdot (l_{j+1} - l_j))I)^{-2}]$$

Since  $g$  is convex, we have that

$$\left. \frac{dg(t)}{dt} \right|_{t=1} \geq g(1) - g(0) = \text{tr}[(B_j - l_{j+1} I)^{-1}] - \text{tr}[(B_j - l_j I)^{-1}].$$

We can conclude that

$$\begin{aligned} \mathbb{E}_{x_j \in D_j} [\text{tr}[(B_{j+1} - l_{j+1} I)^{-1}]] &\leq \text{tr}[(B_j - l_{j+1} I)^{-1}] - \frac{(1 - 1/c)\gamma}{\Phi_j} \text{tr}[(B_j - l_{j+1} I)^{-2}] \\ &\leq \text{tr}[(B_j - l_{j+1} I)^{-1}] \\ &\quad + (1 - 1/c)(1 + \gamma)(\text{tr}[(B_j - l_j I)^{-1}] - \text{tr}[(B_j - l_{j+1} I)^{-1}]) \\ &\leq \text{tr}[(B_j - l_j I)^{-1}] \end{aligned}$$

where the last step follows from  $(1 - 1/c)(1 + \gamma) \geq 1$  and  $\text{tr}[(B_j - l_j I)^{-1}] - \text{tr}[(B_j - l_{j+1} I)^{-1}] \leq 0$ .  $\square$

**Lemma E.7.** *There exists a constant  $C$  such that with probability at least 0.99, Procedure RANDOMIZEDSAMPLING takes at most  $k = C \cdot \bar{d}/\gamma^2$  random points  $x_1, \dots, x_k$  and guarantees that  $\frac{r_k}{l_k} \leq 1 + 8\gamma$ .*

*Proof.*

$$\begin{aligned} \Pr[\text{algorithm finishes within } k \text{ iterations}] &\geq \Pr[4\bar{d}/\gamma + \sum_{j=0}^{k-1} (\frac{\gamma}{\Phi_j(1 - \gamma)} - \frac{\gamma}{\Phi_j(1 + \gamma)}) \geq 8\bar{d}/\gamma] \\ &\geq \Pr[\sum_{j=0}^{k-1} \frac{1}{\Phi_j} \geq \frac{4\bar{d}(1 - \gamma^2)}{\gamma^3}] \\ &\geq \Pr[\frac{\gamma^3 k^2}{4\bar{d}(1 - \gamma^2)} \geq \sum_{j=0}^{k-1} \Phi_j] \end{aligned}$$

By Lemma E.5, every picked matrix  $w_j w_j^\top$  in iteration  $j$  satisfies

$$0 \preceq w_j w_j^\top \preceq \frac{1}{2} \cdot (r_j I - A) \quad (8)$$

with probability at least  $(1 - \gamma')^k \geq 1 - k\gamma'$ . Under the condition of Eq. (8), by Lemma E.6, we have that

$$\mathbb{E}[\sum_{j=0}^{k-1} \Phi_j] \leq k\gamma.$$

Therefore, it holds that

$$\begin{aligned} & \Pr[\text{algorithm does not finish within } k \text{ iterations}] \\ & \leq \Pr[\frac{\gamma^3 k^2}{4\bar{d}(1 - \gamma^2)} \leq \sum_{j=0}^{k-1} \Phi_j] \\ & \leq \Pr[\sum_{j=0}^{k-1} \Phi_j \geq \frac{\gamma^3 k^2}{4\bar{d}(1 - \gamma^2)} \text{ and } \forall j : w_j w_j^\top \preceq \frac{1}{2}(r_j I - A_j)] + \Pr[\exists j : w_j w_j^\top \not\preceq \frac{1}{2}(r_j I - A_j)] \\ & \leq \frac{4\bar{d}(1 - \gamma^2)}{\gamma^2 k} + k\gamma' \\ & \leq \frac{1}{200} \end{aligned}$$

where the last step follows from  $\gamma' \leq 1/(400k)$  and  $k \geq 1600\bar{d}(1 - \gamma^2)/\gamma^2$ .

Let  $\Delta_{r,j} = r_j - r_{j-1}$ ,  $\Delta_{l,j} = l_j - l_{j-1}$ . Since

$$\frac{\Delta_{r,j+1} - \Delta_{l,j+1}}{\Delta_{r,j+1}} = \frac{\gamma/(\Phi_j(1 - \gamma)) - \gamma/(\Phi_j(1 + \gamma))}{\gamma/(\Phi_j(1 - \gamma))} = \frac{2\gamma}{1 + \gamma}.$$

We can claim that

$$\begin{aligned} \frac{r_k - l_k}{r_k} &= \frac{4\bar{d}/\gamma + \sum_{j=0}^{k-1} (\Delta_{r,j+1} - \Delta_{l,j+1})}{2\bar{d}/\gamma + \sum_{j=0}^{k-1} \Delta_{r,j+1}} \\ &\leq \frac{4\bar{d}/\gamma + \sum_{j=0}^{k-1} (\Delta_{r,j+1} - \Delta_{l,j+1})}{2\bar{d}/\gamma + (1 + \gamma)/(2\gamma) \cdot \sum_{j=0}^{k-1} (\Delta_{r,j+1} - \Delta_{l,j+1})} \end{aligned}$$

By the ending condition of the algorithm, it holds that  $r_k - l_k \geq 8\bar{d}/\gamma$ . As a result

$$\frac{r_k - l_k}{r_k} \leq 4\gamma,$$

and

$$\frac{r_k}{l_k} \leq \frac{1}{1 - 4\gamma} \leq 1 + 8\gamma.$$

where the last step follows from  $\gamma \leq 1/8$ . □

## F PERFORMANCE OF I.I.D. DISTRIBUTIONS

**Lemma F.1.** *Let  $D'$  be an arbitrary distribution over  $G$ . There exists an absolute constant  $C$  such that for any  $n \in \mathbb{N}^+$ ,  $\mathcal{F}_{\text{nn}}$  of dimension  $\bar{d}$ ,  $\varepsilon \in (0, 1)$  and  $\delta \in (0, 1)$ , when  $S = (x_1, \dots, x_k)$  are independently from the distribution  $D'$  with  $k \geq \frac{C}{\varepsilon^2} \cdot K_{D'} \log \frac{\bar{d}}{\delta}$  and  $u_j = \frac{D(x_j)}{k \cdot D'(x_j)}$  for each  $j \in [k]$ , the  $k \times \bar{d}$  matrix  $A(i, j) = \sqrt{u_i} \cdot v_j(x_i)$  satisfies*

$$\|A^* A - I\| \leq \varepsilon \text{ with probability at least } 1 - \delta.$$

*Proof.* At the same time, for any fixed  $x$ ,

$$\sum_{i \in [\bar{d}]} |v_i^{(D')}(x)|^2 = \sup_{\alpha(h)} \frac{|\sum_{i=1}^{\bar{d}} \alpha(h)_i \cdot v_i^{(D')}(x)|^2}{\|\alpha(h)\|_2^2} = \sup_{h \in \mathcal{F}_{nn}} \frac{|h^{(D')}(x)|^2}{\|\alpha(h)\|_2^2}$$

by the tightness of the Cauchy Schwartz inequality. Thus

$$K_{\alpha, D'} \stackrel{\text{def}}{=} \sup_{x \in G} \left\{ \sup_{h \in \mathcal{F}_{nn}: h \neq 0} \frac{|h^{(D')}(x)|^2}{\|\alpha(h)\|_2^2} \right\} \quad \text{indicates} \quad \sup_{x \in G} \sum_{i \in [\bar{d}]} |v_i^{(D')}(x)|^2 = K_{\alpha, D'}. \quad (9)$$

For each point  $x_j$  in  $S$  with weight  $u_j = \frac{D(x_j)}{k \cdot D'(x_j)}$ , let  $A_j$  denote the  $j$ th row of the matrix  $A$ . It is a vector in  $\mathbb{R}^{\bar{d}}$  defined by  $A_j(i) = A(j, i) = \sqrt{u_j} \cdot v_i(x_j) = \frac{v_i^{(D')}(x_j)}{\sqrt{k}}$ . So  $A^* A = \sum_{j=1}^k A_j^* \cdot A_j$ .

For  $A_j^* \cdot A_j$ , it is always  $\succeq 0$ . Notice that the only non-zero eigenvalue of  $A_j^* \cdot A_j$  is

$$\lambda(A_j^* \cdot A_j) = A_j \cdot A_j^* = \frac{1}{k} \left( \sum_{i \in [\bar{d}]} |v_i^{(D')}(x_j)|^2 \right) \leq \frac{K_{\alpha, D'}}{k}$$

from (9).

At the same time, because the expectation of the entry  $(i, i')$  in  $A_j^* \cdot A_j$  is

$$\begin{aligned} \mathbb{E}_{x_j \sim D'}[A(j, i) \cdot A(j, i')] &= \mathbb{E}_{x_j \sim D'} \left[ \frac{v_i^{(D')}(x_j) \cdot v_{i'}^{(D')}(x_j)}{k} \right] \\ &= \mathbb{E}_{x_j \sim D'} \left[ \frac{D(x) \cdot v_i(x_j) \cdot v_{i'}(x_j)}{k \cdot D'(x_j)} \right] \\ &= \mathbb{E}_{x_j \sim D} \left[ \frac{v_i(x_j) \cdot v_{i'}(x_j)}{k} \right]. \end{aligned}$$

We can claim that

$$\begin{aligned} \mathbb{E}_{x_j \sim D'}[A(j, i) \cdot A(j, i')] &= 1/k, \forall i = i', \\ \mathbb{E}_{x_j \sim D'}[A(j, i) \cdot A(j, i')] &\leq \rho/k, \forall i \neq i'. \end{aligned}$$

As a result, we have that

$$\begin{aligned} \lambda_{\min} \left( \sum_{j=1}^k \mathbb{E}[A_j^* \cdot A_j] \right) &\geq 1 - \rho \bar{d}, \\ \lambda_{\max} \left( \sum_{j=1}^k \mathbb{E}[A_j^* \cdot A_j] \right) &\leq 1 + \rho \bar{d}. \end{aligned}$$

Now we apply Theorem A.1 on  $A^* A = \sum_{j=1}^k (A_j^* \cdot A_j)$ :

$$\begin{aligned} &\Pr [\lambda(A^* A) \notin [(1 - \varepsilon)(1 - \rho \bar{d}), (1 + \varepsilon)(1 + \rho \bar{d})]] \\ &\leq \bar{d} \left( \frac{e^{-\varepsilon}}{(1 - \varepsilon)^{1 - \varepsilon}} \right)^{(1 - \rho \bar{d}) / \frac{K_{\alpha, D'}}{k}} + \bar{d} \left( \frac{e^{-\varepsilon}}{(1 + \varepsilon)^{1 + \varepsilon}} \right)^{(1 - \rho \bar{d}) / \frac{K_{\alpha, D'}}{k}} \\ &\leq 2\bar{d} \cdot \exp \left( - \frac{\varepsilon^2 \cdot \frac{k(1 - \rho \bar{d})}{K_{\alpha, D'}}}{3} \right) \\ &\leq \delta. \end{aligned}$$

where the last step follows from  $k \geq \frac{6K_{\alpha, D'} \log \frac{\bar{d}}{\delta}}{\varepsilon^2 (1 - \rho \bar{d})}$ .

Thus, we complete the proof.  $\square$

**Lemma F.2.** Given any distribution  $D'$  with the same support of  $D$  and any  $\epsilon > 0$ , the random sampling procedure with  $k = \Theta(K_{\alpha, D'} \log \bar{d} + \frac{K_{D'}}{\epsilon})$  i.i.d. random samples from  $D'$  and coefficients  $\beta_i = 1/k, \forall i \in [k]$  is an  $\epsilon$ -importance sampling procedure.

*Proof.* Because the coefficient  $\beta_i = 1/k = O(\epsilon/K_{\alpha, D'})$  and  $\sum_i \beta_i = 1$ , this indicates the second property of importance sampling procedure.

Since  $k = \Theta(K_{\alpha, D'} \log \bar{d})$ , by Lemma F.1, we know all eigenvalues of  $A^* \cdot A$  are in  $[1 - 1/4, 1 + 1/4]$  with probability  $1 - 10^{-3}$ . This indicates the first property of importance sampling procedure.  $\square$

## G RESULTS FOR ACTIVE LEARNING

Previous work (Chen & Price, 2019) only consider linear case, we generalize it into NN.

**Lemma G.1.** Consider any dimension  $\bar{d}$  linear space  $\mathcal{F}_{nn}$  of functions from a domain  $G$  to  $\mathbb{R}$ . Let  $(D, Y)$  be a joint distribution over  $G \times \mathbb{R}$  and  $f = \arg \min_{h \in \mathcal{F}_{nn}} \mathbb{E}_{(x, y) \sim (D, Y)} [|y - h(x)|^2]$ . Let  $K_\alpha = \sup_{h \in \mathcal{F}_{nn}: h \neq 0} \frac{\sup_{x \in G} |h(x)|^2}{\|\alpha(h)\|_2^2}$  and  $P$  be a importance sampling procedure terminating in  $m_p(\epsilon)$  rounds with probability  $1 - 10^{-3}$  for  $\mathcal{F}_{nn}$ , distribution  $D$ , and  $\epsilon$ . For any  $\epsilon \in (0, 1/10)$ , Algorithm 1 takes  $O(K \log(\bar{d}) + K/\epsilon)$  unlabeled samples from  $D$  and requests at most  $m_p(\epsilon/8)$  labels to output  $\tilde{f}$  satisfying

$$\mathbb{E}_{x \sim D} [|\tilde{f}(x) - f(x)|^2] \leq \epsilon \cdot \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] \text{ in expectation.}$$

*Proof.* We still use  $\|f\|_{D'}$  to denote  $\sqrt{\mathbb{E}_{x \sim D'} [|f(x)|^2]}$ . By Lemma F.2 with  $D$  and the property of  $P$ , with probability at least  $1 - 2 \cdot 10^{-3}$ ,

$$\|h\|_{D_0}^2 \in [\frac{3}{4C_r}, \frac{5}{4C_l}] \cdot \|h\|_D^2 \text{ for every } h \in \mathcal{F}_{nn}. \quad (10)$$

We condition on Eq. (10) holds from now on.

Let  $y_i$  denote a random label of  $x_i$  from  $Y(x_i)$  for each  $i \in [k_0]$  including the unlabeled samples in the algorithm and the labeled samples in Step 6 of Algorithm 1. Let  $f'$  be defined as

$$f' = \arg \min_{h \in \mathcal{F}_{nn}} \mathbb{E}_{x_i \sim D_0, y_i \sim Y(x_i)} [|y_i - h(x_i)|^2]. \quad (11)$$

Using Eq. (10) and Lemma F.2, we have

$$\mathbb{E}_{(x_1, y_1), \dots, (x_{k_0}, y_{k_0})} [\|f' - f\|_D^2] \leq \epsilon \cdot \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] \text{ from the proof of Theorem D.3.}$$

In the next a few paragraph, we will show that  $\tilde{f}$  of a good output of  $P$  with distribution  $D_0$  guarantees  $\|\tilde{f} - f'\|_{D_0}^2 \lesssim \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2]$  with high probability.

Using Eq. (10) and the guarantee of Procedure  $P$ , we have

$$\mathbb{E}_P [\|\tilde{f} - f'\|_{D_0}^2] \leq \epsilon \cdot \mathbb{E}_{x \sim D_0} [|y_i - f'(x_i)|^2]$$

from the proof of Theorem D.3.

Next we bound the right hand side  $\mathbb{E}_{x_i \sim D_0} [|y_i - f'(x_i)|^2]$  by  $\mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2]$  over the randomness of  $(x_1, y_1), \dots, (x_{k_0}, y_{k_0})$ :

$$\begin{aligned} & \mathbb{E}_{(x_1, y_1), \dots, (x_{k_0}, y_{k_0})} \left[ \mathbb{E}_{x_i \sim D_0} [|y_i - f'(x_i)|^2] \right] \\ & \leq \mathbb{E}_{(x_1, y_1), \dots, (x_{k_0}, y_{k_0})} \left[ 2 \mathbb{E}_{x_i \sim D_0} [|y_i - f(x_i)|^2] + 2\|f - f'\|_{D_0}^2 \right] \\ & \leq 2 \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] + \frac{3}{C_l} \mathbb{E}_{(x_1, y_1), \dots, (x_{k_0}, y_{k_0})} [\|f - f'\|_D^2]. \end{aligned}$$

where the last step follows from Eq. (10).

Hence

$$\begin{aligned} \mathbb{E}_{(x_1, y_1), \dots, (x_{k_0}, y_{k_0})} [\mathbb{E}_{\tilde{P}} [\|\tilde{f} - f'\|_{D_0}^2]] &\lesssim \varepsilon (2 + \frac{3\varepsilon}{C_l}) \cdot \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] \\ &\lesssim \varepsilon \cdot \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2]. \end{aligned}$$

where the last step follows from  $\rho \leq 1/2$ .

From all discussion above, by rescaling  $\varepsilon$ , we have

$$\begin{aligned} &\|\tilde{f} - f\|_D^2 \\ &\leq 2\|\tilde{f} - f'\|_D^2 + 2\|f' - f\|_D^2 \\ &\leq \frac{8C_r}{3} \|\tilde{f} - f'\|_{D_0}^2 + 2\|f' - f\|_D^2 \\ &\leq \frac{8C_r\varepsilon}{3} \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] + \frac{\varepsilon}{4} \cdot \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] \\ &\lesssim \varepsilon (1 + C_r) \cdot \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] \\ &\lesssim \varepsilon (1 + \rho\bar{d}) \cdot \mathbb{E}_{(x, y) \sim (D, Y)} [|y - f(x)|^2] \end{aligned}$$

□

**Theorem G.2.** Let  $\mathcal{F}_{\text{nn}}$  be a neural network function family of functions from a domain  $G$  to  $\mathbb{R}$  with dimension  $\bar{d}$ , and consider any (unknown) distribution on  $(x, y)$  over  $G \times \mathbb{R}$ . Let  $D$  be the marginal distribution over  $x$ , and suppose it has bounded “condition number”

$$K := \sup_{h \in \mathcal{F}_{\text{nn}}: h \neq 0} \frac{\sup_{x \in G} |h(x)|^2}{\|h\|_D^2}. \quad (12)$$

Let  $f^* \in \mathcal{F}_{\text{nn}}$  minimize  $\mathbb{E}[|f(x) - y|^2]$ . For any  $\varepsilon < 1$ , there exists an randomized algorithm that takes  $O((1 + \rho\bar{d})(K \log(\bar{d}) + K/\varepsilon))$  unlabeled samples from  $D$  and requires  $O(\bar{d}/\varepsilon)$  labels to output  $\tilde{f}$  such that

$$\mathbb{E}_{\tilde{f}} \mathbb{E}_{x \sim D} [|\tilde{f}(x) - f^*(x)|^2] \lesssim \varepsilon (1 + \rho\bar{d}) \cdot \mathbb{E}_{x, y} [|y - f^*(x)|^2].$$

*Proof.* We can claim that

$$\sup_{h \in \mathcal{F}_{\text{nn}}: h \neq 0} \frac{\sup_{x \in G} |h(x)|^2}{\|h\|_D^2} \leq C_r \sup_{h \in \mathcal{F}_{\text{nn}}: h \neq 0} \frac{\sup_{x \in G} |h(x)|^2}{\|\alpha(h)\|_2^2}$$

where  $C_r \leq 1 + \rho\bar{d}$ .

By applying Lemma G.1, Lemma E.3, we can finish our proof immediately. □

**Theorem G.3.** Let  $f_{\text{nn}}(W, x)$  be a neural network as defined in Definition B.1, and consider any (unknown) distribution on  $(x, y)$  over  $\mathbb{R}^d \times \mathbb{R}$ . Let  $D$  be the marginal distribution over  $x$ , and suppose it has bounded “condition number”

$$K := \sup_{W \in \mathbb{R}^{d \times m}: W \neq 0} \frac{\sup_{x \in G} |f(W, x)|^2}{\|f(W, x)\|_D^2}. \quad (13)$$

Let  $W^* \in \mathbb{R}^{d \times m}$  minimize  $\mathbb{E}[|f_{\text{nn}}(W, x) - y|^2]$ . There exists  $\rho \in (0, 1/10)$ ,  $\bar{d} \geq 3$ ,  $\varepsilon_0 \in (0, 1/10)$ . For any  $0 < \varepsilon \lesssim 1/\log^3(\bar{d})$ , there exists an randomized algorithm  $P$  that takes  $O((1 + \rho\bar{d})(K \log(\bar{d}) + K/\varepsilon))$  unlabeled samples from  $D$  and requires  $O(\bar{d}/\varepsilon)$  labels to output  $\tilde{W} \in \mathbb{R}^{d \times m}$  such that

$$\mathbb{E}_P \mathbb{E}_{x \sim D} [|f_{\text{nn}}(\tilde{W}, x) - f_{\text{nn}}(W^*, x)|^2] \lesssim \varepsilon + \varepsilon \cdot \mathbb{E}_{x, y} [|y - f_{\text{nn}}(W^*, x)|^2].$$

*Proof.* By Claim B.4, we have that there exists  $\{v_1, \dots, v_{\bar{d}}\}$  forms a fixed  $\rho$ -nearly orthonormal basis of  $\mathcal{F}_{nn}$ , such that there exists  $h^*$  satisfied that

$$\|h^*(x) - f_{nn}(W^*, x)\|_D^2 \leq \varepsilon_0.$$

By Theorem G.2, there exists an randomized algorithm that takes  $O((1 + \rho\bar{d})(K \log(\bar{d}) + K/\varepsilon))$  unlabeled samples from  $D$  and requires  $O(\bar{d}/\varepsilon)$  labels to output  $\tilde{h}$  such that

$$\mathbb{E}_{\tilde{h}} \mathbb{E}_{x \sim D} [|\tilde{h}(x) - h^*(x)|^2] \lesssim \varepsilon \cdot \mathbb{E}_{x,y} [|y - h^*(x)|^2].$$

By Claim B.4, there exists  $\tilde{W}$  such that

$$\|\tilde{h}(x) - f_{nn}(\tilde{W}, x)\|_D^2 \leq \varepsilon_0.$$

As a result,

$$\begin{aligned} & \mathbb{E}_P \mathbb{E}_{x \sim D} [|f_{nn}(\tilde{W}, x) - f_{nn}(W^*, x)|^2] \\ & \leq \mathbb{E}_P \mathbb{E}_{x \sim D} [2|h^*(x) - f_{nn}(W^*, x)|^2 + 2|f_{nn}(\tilde{W}, x) - h^*(x)|^2] \\ & \leq \mathbb{E}_P \mathbb{E}_{x \sim D} [2|h^*(x) - f_{nn}(W^*, x)|^2 + 4|\tilde{h}(x) - h^*(x)|^2 + 4|\tilde{h}(x) - f_{nn}(\tilde{W}, x)|^2] \\ & \lesssim \varepsilon_0 + \varepsilon \cdot \mathbb{E}_{x,y} [|y - h^*(x)|^2] \\ & \leq \varepsilon_0 + \varepsilon \cdot \mathbb{E}_{x,y} [2|y - f_{nn}(W^*, x)|^2 + 2|f_{nn}(W^*, x) - h^*(x)|^2] \\ & \lesssim \varepsilon_0 + \varepsilon \cdot \mathbb{E}_{x,y} [|y - f_{nn}(W^*, x)|^2]. \end{aligned}$$

□