

A Additional experiments

Figure 5 and 6 plot the same phenomena as Figure 2 and 3. The only difference is that they strictly follows the setting described in Section 2 where the training data and test data has been normalized to $\|x\| = 1$, and there's no bias term when initializing the neural network linear layer to consist of our problem setup. For each experimental setting (i.e, for each value of d , C_0 , and m), we report results averaged over 5 independent random runs.

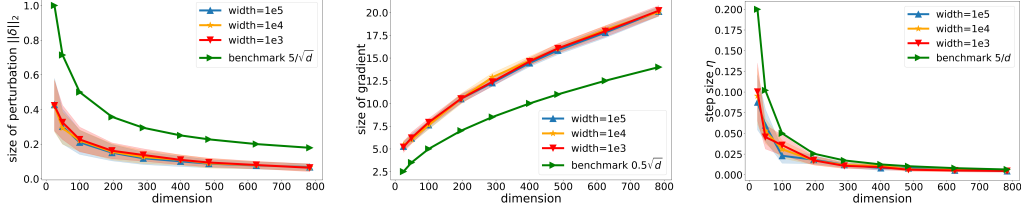


Figure 5: Normalize data $\|x\| = 1$. Smallest size of perturbation to switch the prediction (**left**), norm of the gradient after training the neural networks (**middle**), smallest step size (**right**), as a function of input dimension d for fix $C_0 = 10$ with different width m .

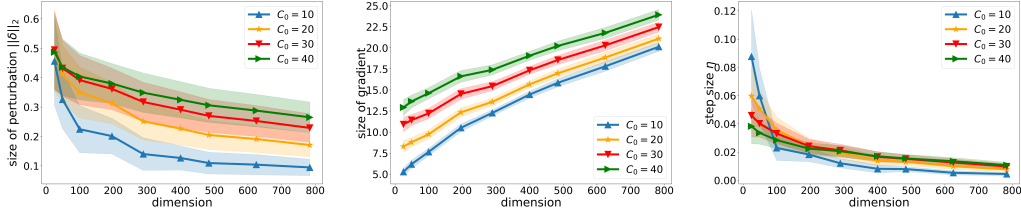


Figure 6: Normalize data $\|x\| = 1$. Smallest size of perturbation to switch the prediction (**left**), norm of the gradient after training the neural networks (**middle**), smallest step size (**right**), as a function of input dimension d for fix $m = 10^5$ with different C_0 .

Figure 7 is the histogram of smallest size of perturbation to switch the prediction, the norm of the gradient after training the neural networks, and the step size when $d = 784$, $m = 10^5$, $C_0 = 10$. The histogram exhibits a Gaussian distribution, and the step size to flip the prediction is small. (η concentrates around 0.045.)

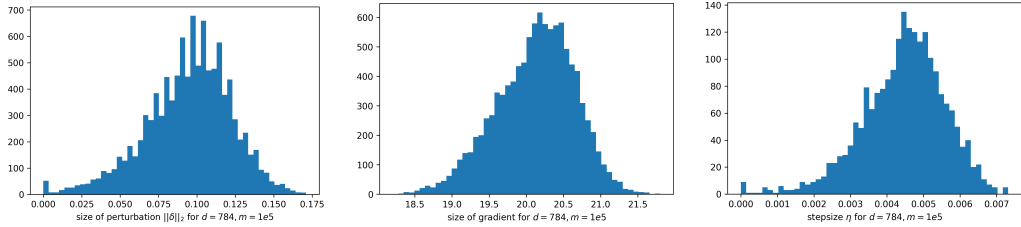


Figure 7: Histogram of smallest size of perturbation to switch the prediction, the norm of the gradient after training the neural networks, and the step size when $d = 784$, $m = 10^5$, $C_0 = 10$.

In Algorithm 1, we describe our projected adversarial training algorithm in details. For generating adversarial example with budget R at each round, we choose learning rate $\alpha = 2.5 \times R/100$ with $T_2 = 100$. This is a common choice, which is first introduced in [Madry et al., 2018]. For updating the weight matrix, we choose learning rate $\beta = 0.01$. We stop when the robust training accuracy is not increasing.

Algorithm 1 Projected Adversarial Training

```
1: Training samples  $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_i, y_i\}_{i=1}^n$ . Initialize  $\mathbf{w}_{s,0} \sim N(0, \mathbf{I}_d)$ ,  $a_s \sim \text{unif}\{\pm 1\}$ ,  $\forall s \in [m]$ 
   with fixed  $\mathbf{a}$ . Epochs  $T_1, T_2$ . Learning rate  $\alpha, \beta$ , perturbation budget per sample  $R$ . Logistic loss
    $\ell$ . Batch size  $bs$ .  $\mathcal{P}_\Delta$  is the projection operator onto the set  $\Delta$ .
2: for  $t = 1, \dots, T_1$  do
3:   for  $k = 1, \dots, T_2$  do
4:      $\tilde{\mathbf{X}} = \mathbf{X} + \alpha \sum_{i=1}^n \nabla_{\mathbf{x}_i} \ell(y_i f(\mathbf{x}_i; \mathbf{a}, \mathbf{W}_{t-1}))$ 
5:      $\tilde{\mathbf{X}} = \mathcal{P}_{\mathcal{B}_{2,\infty}(\mathbf{X}, R)}(\tilde{\mathbf{X}})$  {generate adversarial examples}
6:   end for
7:   for  $\lfloor \frac{n}{bs} \rfloor$  rounds do
8:     Sample a mini-batch of size  $bs$  from  $(\tilde{\mathbf{X}}, \mathbf{y})$  as  $(\tilde{\mathbf{x}}_{ij}, y_{ij})_{j=1}^{bs}$ .
9:      $\mathbf{W}_t = \mathbf{W}_{t-1} - \beta \nabla_{\mathbf{W}_{t-1}} \sum_{j=1}^{bs} \ell(y_{ij} f(\tilde{\mathbf{x}}_{ij}; \mathbf{a}, \mathbf{W}_{t-1}))$ 
10:   end for
11:    $\mathbf{W}_t = \mathcal{P}_{\mathcal{B}_{2,\infty}(\mathbf{W}_0, \frac{C_0}{\sqrt{m}})}(\mathbf{W}_t)$  {Project the weight matrix to satisfy lazy regime.}
12: end for
13: return:  $\mathbf{W}_{T_1}$ 
```

B Proof of theorems

Proof of Theorem 3.1. From Lemma B.5, B.7, B.9, we have that

$$\begin{aligned} \|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| &\geq C'_1 \sqrt{d}, \\ |f(\mathbf{x}; \mathbf{a}, \mathbf{W})| &\leq C'_2, \\ \|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \nabla f(\mathbf{x} + \delta; \mathbf{a}, \mathbf{W})\| &\leq C'_3 \sqrt{d} \text{ for } \|\delta\| \leq o\left(\frac{1}{\sqrt{d}}\right) \end{aligned}$$

We simplify the notation as $f(\mathbf{x}) = f(\mathbf{x}; \mathbf{a}, \mathbf{W})$. Define $\tilde{\eta} = \eta \|\nabla f(\mathbf{x})\|^2$. Without loss of generality, assume $f(\mathbf{x}) > 0$, let

$$\tilde{\eta} = -\frac{2f(\mathbf{x})}{1 - \sup_{\|\delta\| \leq \frac{\tilde{\eta}}{\|\nabla f(\mathbf{x})\|}} \frac{\|\nabla f(\mathbf{x} + \delta)\| - \|f(\mathbf{x})\|}{\|\nabla f(\mathbf{x})\|}},$$

we have for the point $\mathbf{x} + \tilde{\eta} \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|^2}$,

$$\begin{aligned} &f\left(\mathbf{x} + \tilde{\eta} \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|^2}\right) \\ &= f(\mathbf{x}) + \int_0^1 f\left(\mathbf{x} + t\tilde{\eta} \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|^2}\right)' dt \quad (\text{Fundamental theorem of calculus}) \\ &= f(\mathbf{x}) + \int_0^1 \tilde{\eta} \frac{\nabla f(\mathbf{x})^\top}{\|\nabla f(\mathbf{x})\|^2} \nabla f\left(\mathbf{x} + t\tilde{\eta} \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|^2}\right) dt \\ &= f(\mathbf{x}) + \tilde{\eta} + \int_0^1 \tilde{\eta} \frac{\nabla f(\mathbf{x})^\top}{\|\nabla f(\mathbf{x})\|} \frac{(\nabla f(\mathbf{x} + t\tilde{\eta} \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|^2}) - \nabla f(\mathbf{x}))}{\|\nabla f(\mathbf{x})\|} dt \\ &\leq f(\mathbf{x}) + \tilde{\eta} + |\tilde{\eta}| \int_0^1 \frac{\|\nabla f(\mathbf{x} + t\tilde{\eta} \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|^2}) - \nabla f(\mathbf{x})\|}{\|\nabla f(\mathbf{x})\|} dt \quad (\text{Cauchy-Schwarz}) \\ &\leq f(\mathbf{x}) + \tilde{\eta} + |\tilde{\eta}| \sup_{\delta \in \mathbb{R}^d: \|\delta\| \leq \frac{\tilde{\eta}}{\|\nabla f(\mathbf{x})\|}} \frac{\|\nabla f(\mathbf{x} + \delta) - \nabla f(\mathbf{x})\|}{\|\nabla f(\mathbf{x})\|} \\ &= -f(\mathbf{x}) < 0 \end{aligned}$$

Note $|\tilde{\eta}| = \Theta(1)$ and thus $\eta \leq O(\frac{1}{d})$, $\|\eta \nabla f(\mathbf{x})\| = \|\tilde{\eta} \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|^2}\| \leq O(\frac{1}{\sqrt{d}})$.

□

In order to proof the main theorem, we start by giving some standard theorems, which will be used in later proof.

Theorem B.1 (Bernstein's inequality). Let z_1, \dots, z_n be independent real-valued random variables. Assume that there exist positive numbers v and c such that

$$\sum_{i=1}^n \mathbb{E}[z_i^2] \leq v \text{ and } \sum_{i=1}^n \mathbb{E}[|z_i|^q] \leq \frac{q!}{2} v c^{q-2} \text{ for all integers } q \geq 3.$$

If $S = \sum_{i=1}^n (z_i - \mathbb{E}z_i)$, then for all $t > 0$,

$$\mathbb{P}(S \geq \sqrt{2vt} + ct) \leq e^{-t}.$$

We will also use repeatedly that

$$\mathbb{E}_{z \sim \mathcal{N}(0,1)}[|z|^q] \leq (q-1)!! \leq \frac{q!}{2}, \quad (1)$$

as well as the following concentration of χ^2 random variables: let z_1, \dots, z_m be i.i.d. standard Gaussians, then with probability at least $1 - \gamma$, one has

$$\left| \sum_{s=1}^m z_s^2 - m \right| \leq 4\sqrt{m \log(2/\gamma)} \quad (2)$$

Theorem B.2. (Chernoff's inequality) If z_1, z_2, \dots, z_N are independent Bernoulli random variables with parameters μ_i . Let $S_N = \sum_i z_i$ and $p = \mathbb{E}S_N = \sum_i \mu_i$, then for $t > p$,

$$\mathbb{P}(S_N \geq t) \leq \exp(-p) \left(\frac{ep}{t} \right)^t$$

Theorem B.3. (Theorem 2.26 in [Wainwright \[2019\]](#)) Let $\mathbf{z} = (z_1, z_2, \dots, z_n)$ be a vector of i.i.d standard Gaussian variables, and let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be L -Lipschitz w.r.t. the Euclidean norm. Then we have

$$\mathbb{P}(|f(\mathbf{z}) - \mathbb{E}f(\mathbf{z})| \geq t) \leq 2e^{-\frac{t^2}{2L^2}} \text{ for all } t \geq 0$$

In Lemma B.4, we let S_v denote the set of neurons that change sign between weights \mathbf{w}_0 and \mathbf{w} for the \mathbf{x} , S'_v as the set of neurons that change sign between weights \mathbf{w}_0 and \mathbf{w} for the $\mathbf{x} + \delta$. We show that the size of S_v and S'_v is small as long as the weights stay close to initialization.

Lemma B.4. Define the following

$$\begin{aligned} S_v &:= \{s | \exists \mathbf{W}, \mathbf{W} \in \mathcal{B}_{2,\infty}(\mathbf{W}_0, V), \mathbb{1}[\langle \mathbf{w}_s, \mathbf{x} \rangle > 0] \neq \mathbb{1}[\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle > 0]\} \\ S'_v &:= \{s | \exists \delta, \|\delta\| \leq R \leq 0.5, \exists \mathbf{W}, \mathbf{W} \in \mathcal{B}_{2,\infty}(\mathbf{W}_0, V), \mathbb{1}[\langle \mathbf{w}_s, \mathbf{x} + \delta \rangle > 0] \neq \mathbb{1}[\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle > 0]\} \end{aligned}$$

Then with probability at least $1 - \gamma$, the following hold:

$$\begin{aligned} |S_v| &\leq \left| \{s | |\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V\} \right| = \sum_{s=1}^m \mathbb{1}[|\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V] \leq Vm + \sqrt{\frac{m \log(1/\gamma)}{2}} \\ |S'_v| &\leq \left| \{s | |\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V\} \right| = \sum_{s=1}^m \mathbb{1}[|\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V] \leq \left(Vm + \sqrt{\frac{m \log(1/\gamma)}{2}} \right) (1 + R) \end{aligned}$$

Proof of Lemma B.4. Define $\mathbf{v}_s = \mathbf{w}_s - \mathbf{w}_{s,0}$. Note that $s \in S_v$ implies $\forall 1 \leq s \leq m$,

$$|\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq \sup_{\|\mathbf{v}_s\| \leq V} |\langle \mathbf{w}_s - \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq \sup_{\|\mathbf{v}_s\| \leq V} \|\mathbf{v}_s\|_2 \|\mathbf{x}\|_2 \leq V$$

Thus we have

$$\mathbb{E} \left[\frac{1}{m} \sum_{s=1}^m \mathbb{1}[|\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V] \right] = \mathbb{P}(|\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V) \leq \frac{2V}{\|\mathbf{x}\|_2^2 \sqrt{2\pi}} \leq V,$$

where the expectation is with respect to the randomness in initialization, and the inequality holds since $\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle$ is a Gaussian r.v. with variance 1. By Hoeffding inequality, with probability at least $1 - \gamma$,

$$\begin{aligned} & \frac{1}{m} \sum_{s=1}^m \mathbb{1} [|\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V] \\ & \leq \mathbb{E} \left[\frac{1}{m} \sum_{s=1}^m \mathbb{1} [|\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V] \right] + \sqrt{\frac{\log(1/\gamma)}{2m}} \quad (\text{Hoeffding inequality}) \\ & \leq V + \sqrt{\frac{\log(1/\gamma)}{2m}} \end{aligned}$$

As a result, with probability at least $1 - \gamma$, we arrive at the following upperbound on the size of S_v :

$$|S_v| \leq \left| \{s \mid |\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V\} \right| = \sum_{s=1}^m \mathbb{1} [|\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle| \leq V] \leq Vm + \sqrt{\frac{m \log(1/\gamma)}{2}}$$

Same way we can bound the size of S'_v . $s \in S'_v$ implies $\forall 1 \leq s \leq m$,

$$|\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq \sup_{\|\mathbf{v}_s\| \leq V, \|\delta\| \leq R} |\langle \mathbf{v}_s, \mathbf{x} + \delta \rangle| \leq \sup_{\|\mathbf{v}_s\| \leq V, \|\delta\| \leq R} \|\mathbf{v}_s\|_2 \|\mathbf{x} + \delta\|_2 \leq V(1 + R)$$

Thus we have

$$\mathbb{E} \left[\frac{1}{m} \sum_{s=1}^m \mathbb{1} [|\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V(1 + R)] \right] = P(|\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V(1 + R)) \leq V(1 + R),$$

where the expectation is with respect to the randomness in initialization, and the inequality holds since $\langle \mathbf{w}_{s,0}, \mathbf{x} \rangle$ is a Gaussian r.v. with variance 1. By Hoeffding inequality, with probability at least $1 - \gamma$,

$$\begin{aligned} & \frac{1}{m} \sum_{s=1}^m \mathbb{1} [|\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V(1 + R)] \\ & \leq \mathbb{E} \left[\frac{1}{m} \sum_{s=1}^m \mathbb{1} [|\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V(1 + R)] \right] + \sqrt{\frac{\log(1/\gamma)}{2m}} (1 + R) \quad (\text{Hoeffding inequality}) \\ & \leq V(1 + R) + \sqrt{\frac{\log(1/\gamma)}{2m}} (1 + R) \end{aligned}$$

As a result, with probability at least $1 - \gamma$, we arrive at the following upperbound on the size of S_v :

$$|S_v| \leq \left| \{s \mid |\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V\} \right| = \sum_{s=1}^m \mathbb{1} [|\langle \mathbf{w}_{s,0}, \mathbf{x} + \delta \rangle| \leq V] \leq \left(Vm + \sqrt{\frac{m \log(1/\gamma)}{2}} \right) (1 + R)$$

□

Now we begin our proof. Essentially we want to lower bound $\|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\|$, upper bound $|f(\mathbf{x}; \mathbf{a}, \mathbf{W})|$ and upper bound $\|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \nabla f(\mathbf{x} + \delta; \mathbf{a}, \mathbf{W})\|$. Lemma B.5 gives the upper bound of $|f(\mathbf{x}; \mathbf{a}, \mathbf{W})|$ as $O(1)$.

Lemma B.5. For any \mathbf{x} , with probability at least $1 - \gamma$, the following holds for all $\mathbf{W} \in \mathcal{B}_{2,\infty} \left(\mathbf{W}_0, \frac{C_0}{\sqrt{m}} \right)$,

$$|f(\mathbf{x}; \mathbf{a}, \mathbf{W})| \leq \sqrt{2 \log(2/\gamma)} + \frac{2 \log(2/\gamma)}{\sqrt{m}} + C_0$$

Particularly, there exists $C > 0$ such that for $m \geq C \log(2/\gamma)$, we have

$$|f(\mathbf{x}; \mathbf{a}, \mathbf{W})| \leq 2\sqrt{\log(2/\gamma)} + C_0$$

Proof of Lemma B.5. From Bubeck et al. [2021] we know that $|f(x; a, W_0)| \leq \sqrt{2 \log(2/\gamma)} + \frac{2 \log(2/\gamma)}{\sqrt{m}}$. Now we consider bound $|f(x; a, W) - f(x; a, W_0)|$,

$$\begin{aligned}
|f(x; a, W) - f(x; a, W_0)| &\leq \left| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s (\sigma(w_s^\top x) - \sigma(w_{s,0}^\top x)) \right| \\
&\leq \frac{1}{\sqrt{m}} \sum_{s=1}^m |\sigma(w_s^\top x) - \sigma(w_{s,0}^\top x)| && \text{(Triangle Inequality)} \\
&\leq \sqrt{m} |w_s^\top x - w_{s,0}^\top x| && (\sigma(\cdot) \text{ is 1-Lipschitz.}) \\
&\leq \sqrt{m} \|w_s - w_{s,0}\| \|x\| && (\|w_s - w_{s,0}\| \leq \frac{C_0}{\sqrt{m}}) \\
&\leq C_0
\end{aligned}$$

Thus,

$$|f(x; a, W)| = |f(x; a, W_0)| + |f(x; a, W) - f(x; a, W_0)| \leq \sqrt{2 \log(2/\gamma)} + \frac{2 \log(2/\gamma)}{\sqrt{m}} + C_0$$

□

In Lemma B.6, we want to calculate the upper bound of the probability that the neuron flips sign due to (1) perturbation δ ; (2) different perturbation δ, δ' on the data x at weight w_s .

Lemma B.6. For any δ such that $\|\delta\| \leq R \leq \frac{1}{2}$,

$$P\left(\text{sign}(w_{s,0}^\top x) \neq \text{sign}(w_{s,0}^\top (x + \delta))\right) \leq R\sqrt{2 \log(d)} + \frac{1}{d} \quad (3)$$

$$\begin{aligned}
P\left(\exists \delta', \delta' \in \mathcal{B}_2(\delta, \varepsilon), \text{ and } \exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), \text{sign}(w_s^\top (x + \delta)) \neq \text{sign}(w_s^\top (x + \delta'))\right) \\
\leq 2\varepsilon \left(\sqrt{d} + 2\sqrt{d \log(2/\varepsilon)}\right) + (1 + R + \varepsilon)V \quad (4)
\end{aligned}$$

Proof of Lemma B.6. Equation (3) directly follows from Bubeck et al. [2021]. For equation (4), we have

$$\begin{aligned}
&P\left(\exists \delta', \delta' \in \mathcal{B}_2(\delta, \varepsilon), \text{ and } \exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), \text{sign}(w_s^\top (x + \delta)) \neq \text{sign}(w_s^\top (x + \delta'))\right) \\
&\leq P\left(\exists \delta', \delta' \in \mathcal{B}_2(\delta, \varepsilon), \text{ and } \exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), |w_s^\top (\delta' - \delta)| \geq t\right) \\
&\quad + P\left(\exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), |w_s^\top (x + \delta)| \leq t\right) \quad \text{(Holds for any threshold } t.) \\
&\leq P\left(\exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), \|w_s\| \geq t/\varepsilon\right) + P\left(\exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), |w_s^\top (x + \delta)| \leq t\right) \\
&\leq P\left(\exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), \|w_{s,0}\| \geq t/\varepsilon - \|w_s - w_{s,0}\|\right) \\
&\quad + P\left(\exists W, W \in \mathcal{B}_{2,\infty}(W_0, V), |w_{s,0}^\top (x + \delta)| \leq t + |(w_s - w_{s,0})^\top (x + \delta)|\right) \\
&\quad \quad \quad \text{(Triangle Inequality, pick } t = \varepsilon \left(\sqrt{d + 4\sqrt{d \log(2/\varepsilon)}} + V\right)) \\
&\leq P\left(\|w_{s,0}\| \geq \sqrt{d + 4\sqrt{d \log(2/\varepsilon)}}\right) + P\left(|w_{s,0}^\top (x + \delta)| \leq \varepsilon \sqrt{d + 4\sqrt{d \log(2/\varepsilon)}} + (1 + R + \varepsilon)V\right) \\
&\quad \quad \quad (w_{s,0}^\top (x + \delta) \sim \mathcal{N}(0, \sigma^2) \text{ with } \sigma^2 \geq \frac{1}{2} \text{ since } \|\delta\| \leq \frac{1}{2}) \\
&\leq \varepsilon + \varepsilon \sqrt{d + 4\sqrt{d \log(2/\varepsilon)}} + (1 + R + \varepsilon)V \quad \text{(Using equation (2).)} \\
&= 2\varepsilon(\sqrt{d} + 2\sqrt{d \log(2/\varepsilon)}) + (1 + R + \varepsilon)V
\end{aligned}$$

□

Lemma B.7 gives lower bound on $\|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| \geq \Omega(\sqrt{d})$.

Lemma B.7. For any \mathbf{x} , with probability at least $1 - \gamma$, the following holds for all $\mathbf{W} \in \mathcal{B}_{2,\infty}\left(\mathbf{W}_0, \frac{C_0}{\sqrt{m}}\right)$,

$$\begin{aligned} \|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| &\geq \left(\frac{1}{2} - \left(\sqrt{\frac{2\log(4/\gamma)}{m}} + \frac{\log(4/\gamma)}{m}\right)\right)^{1/2} \left(d - 5\sqrt{d\log(8/\gamma)}\right)^{1/2} \\ &\quad - C_0 - \sqrt{\frac{1}{\sqrt{m}}}\left(C_0 + \sqrt{\frac{\log(4/\gamma)}{2}}\right) \left(d + 4\sqrt{d\log(8m/\gamma)}\right)^{1/2} \end{aligned}$$

Particularly, there exists $C > 0$ such that for $m \geq C \log(4/\gamma)$ and $d \geq C \frac{\log(8m/\gamma)}{m}$, we have

$$\|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| \geq \frac{1}{4}\sqrt{d}$$

Proof of Lemma B.7. Follow the process of Bubeck et al. [2021], let $\mathbf{P} = \mathbf{I}_d - \mathbf{x}\mathbf{x}^\top$ be the projection on the orthogonal complement of the span of \mathbf{x} . We have $\|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| \geq \|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\|$. Thus we have,

$$\begin{aligned} \|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| &\geq \|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| \\ &= \|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0) + \mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0)\| \\ &\geq \|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0)\| - \|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0)\| \end{aligned} \tag{5}$$

Since $a_s \mathbf{P}\mathbf{w}_{s,0}$ is distributed as $\mathcal{N}(0, \mathbf{I}_{d-1})$, denote $\mathbf{z} := a_s \mathbf{P}\mathbf{w}_{s,0}$, we have

$$\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0) = \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P}\mathbf{w}_{s,0} \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x}) \stackrel{(d)}{=} \left(\sqrt{\frac{1}{m} \sum_{s=1}^m \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})^2} \right) \mathbf{z} \text{ where } \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_{d-1})$$

where $\stackrel{(d)}{=}$ means equal in distribution.

Using (2), we have with probability at least $1 - \gamma$

$$\|\mathbf{z}\|^2 \geq d - 1 - 4\sqrt{d\log(2/\gamma)} \geq d - 5\sqrt{d\log(2/\gamma)} \quad (d \geq 1, \gamma < 2/e)$$

Apply Bernstein's inequality with $v = m, c = 1$, we have with probability at least $1 - \gamma$,

$$\begin{aligned} \frac{1}{m} \sum_{s=1}^m \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})^2 &\geq \mathbb{E}_{X \sim \mathcal{N}(0,1)}[|\sigma'(X)|^2] - \left(\sqrt{\frac{2\log(1/\gamma)}{m}} + \frac{\log(1/\gamma)}{m} \right) \\ &= \frac{1}{2} - \left(\sqrt{\frac{2\log(1/\gamma)}{m}} + \frac{\log(1/\gamma)}{m} \right) \end{aligned}$$

Thus we can get

$$\|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0)\| \geq \left(\frac{1}{2} - \left(\sqrt{\frac{2\log(2/\gamma)}{m}} + \frac{\log(2/\gamma)}{m}\right)\right)^{1/2} \left(d - 5\sqrt{d\log(4/\gamma)}\right)^{1/2}$$

Now we start upper bound $\|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0)\|$. Note that

$$\begin{aligned}
& \|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0)\| \\
&= \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P}(\mathbf{w}_s \sigma'(\mathbf{w}_s^\top \mathbf{x}) - \mathbf{w}_{s,0} \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\leq \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P}(\mathbf{w}_s \sigma'(\mathbf{w}_s^\top \mathbf{x}) - \mathbf{w}_{s,0} \sigma'(\mathbf{w}_s^\top \mathbf{x})) \right\| + \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\leq \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s (\mathbf{w}_s \sigma'(\mathbf{w}_s^\top \mathbf{x}) - \mathbf{w}_{s,0} \sigma'(\mathbf{w}_s^\top \mathbf{x})) \right\| \\
&\quad + \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\leq \sqrt{m} \|\mathbf{w}_s - \mathbf{w}_{s,0}\| + \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\leq C_0 + \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \tag{6}
\end{aligned}$$

Now we start bounding the second term of equation (6). We have that with probability at least $1 - \gamma$,

$$\begin{aligned}
& \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \\
&\stackrel{(d)}{=} \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \left(\sqrt{\frac{1}{m} \sum_{s=1}^m (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x}))^2} \right) \mathbf{z} \quad (\mathbf{z} := a_s \mathbf{P} \mathbf{w}_{s,0} \sim N(0, \mathbf{I}_{d-1})) \\
&= \left(\sqrt{\frac{1}{m} \sum_{s=1}^m \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x}))^2} \right) \mathbf{z} \\
&= \sqrt{\frac{1}{m}} |S_v| \cdot \mathbf{z} \quad (\text{By the definition of } S_v \text{ in Lemma B.4 with } V = \frac{C_0}{\sqrt{m}}) \\
&\leq \sqrt{\frac{1}{\sqrt{m}}} (C_0 + \sqrt{\frac{\log(1/\gamma)}{2}}) \cdot \mathbf{z} \tag{7}
\end{aligned}$$

Using equation (2), we see that with probability at least $1 - \gamma$, one has for all $1 \leq s \leq m$,

$$\|\mathbf{z}\| = \|a_s \mathbf{P} \mathbf{w}_{s,0}\| \leq \sqrt{d + 4\sqrt{d \log(2m/\gamma)}}$$

Thus follow from equation (6), we have

$$\|\mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \mathbf{P}\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}_0)\| \leq C_0 + \sqrt{\frac{1}{\sqrt{m}}} (C_0 + \sqrt{\frac{\log(2/\gamma)}{2}}) \sqrt{d + 4\sqrt{d \log(4m/\gamma)}}$$

And as a result, with probability at least $1 - \gamma$

$$\begin{aligned} \|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W})\| &\geq \left(\frac{1}{2} - \left(\sqrt{\frac{2\log(4/\gamma)}{m}} + \frac{\log(4/\gamma)}{m} \right) \right)^{1/2} \left(d - 5\sqrt{d\log(8/\gamma)} \right)^{1/2} \\ &\quad - C_0 - \sqrt{\frac{1}{\sqrt{m}}} (C_0 + \sqrt{\frac{\log(4/\gamma)}{2}}) \sqrt{d + 4\sqrt{d\log(8m/\gamma)}} \geq \Omega(\sqrt{d}) \end{aligned}$$

□

Now we start bounding $\|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \nabla f(\mathbf{x} + \delta; \mathbf{a}, \mathbf{W})\|$. In order to prove Lemma B.9, we will leverage the fact that for any $h \in \mathbb{R}^d$, we have $\|h\| = \sup_{r \in \mathbb{S}^{d-1}} r \cdot h$. We first start a lemma with fixed r, δ to prove Lemma B.8, then use the ε -net argument to prove Lemma B.9.

Lemma B.8. Fix $r \in \mathbb{S}^{d-1}$, and $\delta \in \mathbb{R}^d$ such that $\|\delta\| \leq R \leq \frac{C_1}{\sqrt{d}}$, C_1 is a constant. For any \mathbf{x} , with probability at least $1 - \gamma$, the following holds for all $\mathbf{W} \in \mathcal{B}_{2,\infty}(\mathbf{W}_0, \frac{C_0}{\sqrt{m}})$,

$$\begin{aligned} &\frac{1}{\sqrt{m}} \sum_{s=1}^m a_s(\mathbf{w}_s \cdot r) (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta))) \\ &\leq 2\sqrt{\log(4/\gamma)} \left((4R\sqrt{\log(d)})^{1/4} + \sqrt{\frac{\log(4/\gamma)}{m}} \right) + C_0 \\ &\quad + 3\sqrt{\frac{1}{\sqrt{m}}} (C_0 + \sqrt{\frac{\log(4/\gamma)}{2}}) \sqrt{d + 4\sqrt{d\log(8m/\gamma)}} + 5(C_0 + \sqrt{\log(4/\gamma)}) \cdot \sqrt{\log(4m/\gamma)} \end{aligned}$$

Proof of Lemma B.8. Let $X_{s,0} = a_s(\mathbf{w}_{s,0} \cdot r) (\sigma'(\mathbf{w}_{s,0}^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta)))$, $X_s = a_s(\mathbf{w}_s \cdot r) (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)))$. We have $\mathbb{E}[X_{s,0}] = 0$, and equation (1) gives us $\mathbb{E}[|\mathbf{w}_{s,0} \cdot r|^{2q}] \leq \frac{(2q)!}{2}$. Thus, we have

$$\begin{aligned} \mathbb{E}[|X_{s,0}|^q] &\leq \mathbb{E}[|\mathbf{w}_{s,0} \cdot r|^q |\sigma'(\mathbf{w}_{s,0}^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))|^q] \\ &\leq \sqrt{\mathbb{E}[|\mathbf{w}_{s,0} \cdot r|^{2q}] \cdot P(\text{sign}(\mathbf{w}_{s,0}^\top \mathbf{x}) \neq \text{sign}(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta)))} \quad (\text{Using equation (3)}) \\ &\leq \sqrt{\frac{(2q)!}{2}} \cdot \sqrt{R\sqrt{2\log(d)} + \frac{1}{d}} \\ &\leq \frac{q!}{2} 2^q \cdot \sqrt{2R\sqrt{2\log(d)}} \end{aligned}$$

$$\begin{aligned} \mathbb{E}[|X_{s,0}|^2] &\leq \sqrt{\mathbb{E}[(\mathbf{w}_{s,0} \cdot r)^4] \cdot P(\text{sign}(\mathbf{w}_{s,0}^\top \mathbf{x}) \neq \text{sign}(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta)))} \\ &\leq 2\sqrt{R\sqrt{2\log(d)}} \quad (\mathbb{E}_{X \sim \mathcal{N}(0,1)}[|X|^4] \leq 3!! \leq 4) \\ &\leq 4\sqrt{R\sqrt{\log(d)}} \end{aligned}$$

Using Theorem B.1 with $v = 4m\sqrt{R\sqrt{\log(d)}}$, $c = 2$, we have

$$\begin{aligned} \frac{1}{\sqrt{m}} \sum_{s=1}^m X_{s,0} &\leq \sqrt{8\sqrt{R\sqrt{\log(d)}} \log(1/\gamma)} + \frac{2\log(1/\gamma)}{\sqrt{m}} \\ &= 2\sqrt{\log(1/\gamma)} \left((4R\sqrt{\log(d)})^{1/4} + \sqrt{\frac{\log(1/\gamma)}{m}} \right) \end{aligned}$$

Now we start bounding $|\frac{1}{\sqrt{m}} \sum_{s=1}^m (X_s - X_{s,0})|$. In fact, here we no longer fix r, δ and directly bound $|\sup_{r \in \mathbb{S}^{d-1}, \delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m (X_s - X_{s,0})|$.

$$\begin{aligned}
& \left| \sup_{\substack{r \in \mathbb{S}^{d-1}, \\ \delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m (X_s - X_{s,0}) \right| \\
&= \left| \sup_{\substack{r \in \mathbb{S}^{d-1}, \\ \delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \left(\mathbf{w}_s^\top r (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta))) - \mathbf{w}_{s,0}^\top r (\sigma'(\mathbf{w}_{s,0}^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right) \right| \\
&\leq \left| \sup_{\substack{r \in \mathbb{S}^{d-1}, \\ \delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \left(\mathbf{w}_s^\top r (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta))) - \mathbf{w}_{s,0}^\top r (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta))) \right) \right| \\
&\quad + \left| \sup_{\substack{r \in \mathbb{S}^{d-1}, \\ \delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \left(\mathbf{w}_{s,0}^\top r (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta))) - \mathbf{w}_{s,0}^\top r (\sigma'(\mathbf{w}_{s,0}^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right) \right| \\
&\hspace{15em} \text{(Triangle Inequality)} \\
&\leq \frac{1}{\sqrt{m}} \sum_{s=1}^m \|\mathbf{w}_s - \mathbf{w}_{s,0}\| |\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta))| + \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\quad + \left\| \sup_{\substack{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right\| \\
&\hspace{15em} \text{(Cauchy-Schwarz, triangle Inequality)} \\
&\leq C_0 + \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \tag{8} \\
&\quad + \left\| \sup_{\substack{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right\| \tag{9}
\end{aligned}$$

Define $\mathbf{P} = \mathbf{I}_d - \mathbf{x}\mathbf{x}^\top$ as the projection on the orthogonal complement of the span of \mathbf{x} . For the second term of equation (9), we have

$$\begin{aligned}
& \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\leq \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| + \left\| \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{x} \mathbf{x}^\top \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\leq \left\| \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| \\
&\quad + \frac{1}{\sqrt{m}} \sum_{s=1}^m \|\mathbf{x} \mathbf{x}^\top \mathbf{w}_{s,0}\| \cdot |\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})| \\
&\leq \left\| \sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P} \mathbf{w}_{s,0} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top \mathbf{x}) - \sigma'(\mathbf{w}_{s,0}^\top \mathbf{x})) \right\| + \frac{1}{\sqrt{m}} |S_v| \max_{1 \leq s \leq m} \|\mathbf{x} \mathbf{x}^\top \mathbf{w}_{s,0}\| \\
&\hspace{15em} \text{(By the definition of } S_v \text{ in Lemma B.4, Hölder's inequality)} \\
&\leq \sqrt{\frac{1}{\sqrt{m}} (C_0 + \sqrt{\frac{\log(4/\gamma)}{2}})} \sqrt{d + 4\sqrt{d \log(8m/\gamma)}} + (C_0 + \sqrt{\log(4/\gamma)}) \cdot \sqrt{\log(4m/\gamma)} \\
&\hspace{15em} \text{(Equation (7), } \mathbf{x}^\top \mathbf{w}_{s,0} \sim N(0, 1))
\end{aligned}$$

where the last line holds because with probability at least $1 - \gamma$, one has for all $1 \leq s \leq m$, $|\mathbf{x}^\top \mathbf{w}_{s,0}| \leq \sqrt{\log(m/\gamma)}$.

Define $\mathbf{P}' = \mathbf{I}_d - \frac{(\mathbf{x}+\delta)(\mathbf{x}+\delta)^\top}{\|\mathbf{x}+\delta\|^2}$ as the projection on the orthogonal complement of the span of $\mathbf{x} + \delta$. We bound the third term of equation (9) the same way as follows.

$$\begin{aligned}
& \left\| \sup_{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right\| \\
& \leq \left\| \sup_{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P}' \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right\| \\
& \quad + \left\| \sup_{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \frac{(\mathbf{x} + \delta)(\mathbf{x} + \delta)^\top}{\|\mathbf{x} + \delta\|^2} \mathbf{w}_{s,0} (\sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right\| \\
& \leq \left\| \sup_{\substack{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}} \\ \forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \mathbf{P}' \mathbf{w}_{s,0} (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))) \right\| \\
& \quad + \left\| \sup_{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \frac{1}{\sqrt{m}} \sum_{s=1}^m |\sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta))| \left\| \frac{(\mathbf{x} + \delta)(\mathbf{x} + \delta)^\top}{\|\mathbf{x} + \delta\|^2} \mathbf{w}_{s,0} \right\| \right\| \\
& \leq \sqrt{\sup_{\substack{\forall s \in [m], \mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}} \\ \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}}} \frac{1}{m} \sum_{s=1}^m (\sigma'((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top (\mathbf{x} + \delta)) - \sigma'(\mathbf{w}_{s,0}^\top (\mathbf{x} + \delta)))} \sup_{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \|\mathbf{z}(\delta)\|} \\
& \quad \text{(Define } \mathbf{z}(\delta) := a_s \mathbf{P}' \mathbf{w}_{s,0} \sim N(0, \mathbf{I}_{d-1}) \text{)} \\
& \quad + \sup_{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \frac{1}{\sqrt{m}} |S'_v| \left\| \frac{(\mathbf{x} + \delta)(\mathbf{x} + \delta)^\top}{\|\mathbf{x} + \delta\|^2} \mathbf{w}_{s,0} \right\| \quad \text{(By definition of } S'_v \text{ in Lemma B.4)} \\
& \leq \sqrt{\frac{1}{m} |S'_v|} \cdot \sup_{\delta \in \mathbb{R}^d, \|\delta\| \leq \frac{C_1}{\sqrt{d}}} \|\mathbf{z}(\delta)\| \\
& \quad + \frac{1}{\sqrt{m}} |S'_v| \cdot 2 \max_{1 \leq s \leq m} \left| \mathbf{x}^\top \mathbf{w}_{s,0} + \delta^\top \mathbf{w}_{s,0} \right| \quad (\|\mathbf{z}(\delta)\| \leq \|\mathbf{w}_{s,0}\|) \\
& \leq \sqrt{\frac{1+R}{\sqrt{m}}} (C_0 + \sqrt{\frac{\log(4/\gamma)}{2}}) \sqrt{d + 4\sqrt{d \log(8m/\gamma)}} \\
& \quad + (C_0 + \sqrt{\log(4/\gamma)}) (1+R) \cdot 2 \left(\sqrt{\log(4m/\gamma)} + \frac{\sqrt{d + 4\sqrt{d \log(8m/\gamma)}}}{\sqrt{d}} \right) \\
& \quad \quad \quad (R \leq \frac{C_1}{\sqrt{d}}, d > \log(m/\gamma) \log(1/\gamma)) \\
& \leq \sqrt{\frac{2}{\sqrt{m}}} (C_0 + \sqrt{\frac{\log(4/\gamma)}{2}}) \sqrt{d + 4\sqrt{d \log(8m/\gamma)}} + 4(C_0 + \sqrt{\log(4/\gamma)}) \left(\sqrt{\log(4m/\gamma)} + 3 \right)
\end{aligned}$$

As a result, we get

$$\begin{aligned}
\frac{1}{\sqrt{m}} \sum_{s=1}^m X_s & \leq \frac{1}{\sqrt{m}} \sum_{s=1}^m X_{s,0} + \left| \frac{1}{\sqrt{m}} \sum_{s=1}^m (X_s - X_{s,0}) \right| \\
& \leq 2\sqrt{\log(4/\gamma)} \left((4R\sqrt{\log(d)})^{1/4} + \sqrt{\frac{\log(4/\gamma)}{m}} \right) + C_0 \\
& \quad + 3\sqrt{\frac{1}{\sqrt{m}}} (C_0 + \sqrt{\frac{\log(4/\gamma)}{2}}) \sqrt{d + 4\sqrt{d \log(8m/\gamma)}}
\end{aligned}$$

$$+ 5(C_0 + \sqrt{\log(4/\gamma)}) \left(\sqrt{\log(4m/\gamma)} + 3 \right)$$

□

Lemma B.9. Let $\|\delta\| \leq R \leq \frac{C_1}{\sqrt{d}}$, $\bar{C} = \max\{1, C_0\}$, $m \geq d^{2.4}$. Then, for any \mathbf{x} , with probability at least $1 - \gamma$, the following holds for all $\mathbf{W} \in \mathcal{B}_{2,\infty} \left(\mathbf{W}_0, \frac{C_0}{\sqrt{m}} \right)$,

$$\sup_{\substack{r \in \mathbb{S}^{d-1}, \\ \delta \in \mathbb{R}^d: \|\delta\| \leq R}} \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s(\mathbf{w}_s \cdot r) (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta))) \leq 9 \left(C_1 d^2 \log^2(md) \sqrt{\frac{\log(d)}{d}} \right)^{1/4} + \frac{15d \log(md)}{\sqrt{m}} + 327 \bar{C} C_0 d^{0.25}$$

provided that $d \geq \log(4/\gamma)^2$. Particularly, for any $c > 0$, there exists C, C' such that if $\log^2(md) \sqrt{\frac{\log(d)}{d}} \leq C$, $\frac{d \log(md)}{\sqrt{m}} \leq C' \sqrt{d}$, then we have,

$$\sup_{\substack{r \in \mathbb{S}^{d-1}, \\ \delta \in \mathbb{R}^d, \|\delta\| \leq R}} \|\nabla f(\mathbf{x}; \mathbf{a}, \mathbf{W}) - \nabla f(\mathbf{x} + \delta; \mathbf{a}, \mathbf{W})\| \leq c\sqrt{d}$$

The above can realize when $m \leq O(\exp(d^{0.24}))$

Proof. Define $\Phi(r, \delta) = \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s(\mathbf{w}_s \cdot r) (\sigma'(\mathbf{w}_s^\top \mathbf{x}) - \sigma'(\mathbf{w}_s^\top (\mathbf{x} + \delta)))$, and N_ε an ε -net for $\Omega = \{(r, \delta), \|r\| = 1, \|\delta\| \leq R\}$. In Lemma B.8 we bound $\Phi(r, \delta)$ for a fixed r and δ , here we bounded it uniformly over Ω . We know that $|N_\varepsilon| \leq (10/\varepsilon)^{2d}$. Using Lemma B.8, we obtain with probability at least $1 - \gamma$,

$$\begin{aligned} \sup_{(r, \delta) \in \Omega} \Phi(r, \delta) &\leq \sup_{(r, \delta) \in N_\varepsilon} \Phi(r, \delta) + \sup_{(r, \delta), (r', \delta') \in \Omega: \|r - r'\| + \|\delta - \delta'\| \leq \varepsilon} |\Phi(r, \delta) - \Phi(r', \delta')| \\ &\leq 2\sqrt{2d \log(10/\varepsilon) + \log(4/\gamma)} \left((4R\sqrt{\log(d)})^{1/4} + \sqrt{\frac{2d \log(10/\varepsilon) + \log(4/\gamma)}{m}} \right) \\ &\quad + C_0 + 3\sqrt{\frac{1}{\sqrt{m}} (C_0 + \sqrt{\frac{\log(4/\gamma)}{2}})} \sqrt{d + 4\sqrt{d \log(8m/\gamma)}} \\ &\quad + 5(C_0 + \sqrt{\log(4/\gamma)}) \cdot \left(\sqrt{\log(4m/\gamma)} + 3 \right) + \sup_{\substack{(r, \delta), (r', \delta') \in \Omega: \\ \|r - r'\| + \|\delta - \delta'\| \leq \varepsilon}} |\Phi(r, \delta) - \Phi(r', \delta')| \end{aligned} \tag{10}$$

For r, r' , one has

$$|\Phi(r, \delta) - \Phi(r', \delta)| \leq \frac{\|r - r'\|}{\sqrt{m}} \sum_{s=1}^m \|\mathbf{w}_s\| \leq \frac{\|r - r'\|}{\sqrt{m}} \sum_{s=1}^m (\|\mathbf{w}_{s,0}\| + \frac{C_0}{\sqrt{m}})$$

Using (2), we know with probability at least $1 - \gamma$, one has for all $s \in [m]$,

$$\|\mathbf{w}_{s,0}\|^2 \leq d + 4\sqrt{d \log(m/\gamma)},$$

so that in this event we have,

$$|\Phi(r, \delta) - \Phi(r', \delta)| \leq \|r - r'\| \left(\sqrt{md + 4m\sqrt{d \log(m/\gamma)}} + C_0 \right)$$

On the other hand, for δ, δ' , we write

$$\begin{aligned} &|\Phi(r, \delta) - \Phi(r, \delta')| \\ &\leq \frac{1}{\sqrt{m}} \left| \sum_{s=1}^m \mathbb{1}\{\text{sign}(\mathbf{w}_s^\top (\mathbf{x} + \delta)) > \text{sign}(\mathbf{w}_s^\top (\mathbf{x} + \delta'))\} a_s \mathbf{w}_s \cdot r \right| \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{\sqrt{m}} \left| \sum_{s=1}^m \mathbb{1}\{\text{sign}(\mathbf{w}_s^\top(\mathbf{x} + \delta)) < \text{sign}(\mathbf{w}_s^\top(\mathbf{x} + \delta'))\} a_s \mathbf{w}_s \cdot \mathbf{r} \right| \\
& \leq \frac{1}{\sqrt{m}} \left\| \sup_{\substack{\forall s \in [m], \\ \mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq V}} \sum_{s=1}^m \mathbb{1}\{\text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top(\mathbf{x} + \delta)) > \text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top(\mathbf{x} + \delta'))\} a_s \mathbf{w}_s \right\| \\
& + \frac{1}{\sqrt{m}} \left\| \sup_{\substack{\forall s \in [m], \\ \mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq V}} \sum_{s=1}^m \mathbb{1}\{\text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top(\mathbf{x} + \delta)) < \text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top(\mathbf{x} + \delta'))\} a_s \mathbf{w}_s \right\|
\end{aligned}$$

Letting $X_s(\delta) = \mathbb{1}\{\exists \delta' : \|\delta - \delta'\| \leq \varepsilon \text{ and } \exists \mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq V, \text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top(\mathbf{x} + \delta)) \neq \text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top(\mathbf{x} + \delta'))\}$, we now control with exponentially high probability $\sum_{s=1}^m X_s(\delta)$. By (4) in Lemma B.6, we know that $X_s(\delta)$ is a Bernoulli of parameter at most $p = 2\varepsilon(\sqrt{d} + 2\sqrt{d\log(2/\varepsilon)}) + (1 + R + \varepsilon)V$. Using Theorem B.2 and apply a union bound, we have

$$P\left(\exists(r, \delta) \in N_\varepsilon : \sum_{s=1}^m X_s(\delta) \geq k\right) \leq \left(\frac{10}{\varepsilon}\right)^{2d} \exp(-pm) \left(\frac{epm}{k}\right)^k := g(p)$$

We would like $g(p) \leq \gamma$. Since $\sqrt{\varepsilon}(1 + 2\sqrt{\log(2/\varepsilon)}) \leq 4\varepsilon^{3/8}$ holds for $\varepsilon \in (0, 1)$, we keep upper bound $p \leq 8\sqrt{d}\varepsilon^{7/8} + 3V$. Note that $p \leq \frac{k}{m}$ to make sure $g(p)$ increases as p increases. Choose $\varepsilon = m^{-4/7}d^{-4/7}$, $V = \frac{C_0}{\sqrt{m}}$, $\bar{C} = \max\{C_0, 1\}$, $k = 44\bar{C}\sqrt{m}$, with $m \geq 58$, $m \geq d^{2.4}$, we have

$$\begin{aligned}
\log g(8\sqrt{d}\varepsilon^{7/8} + \frac{3C_0}{\sqrt{m}}) &= -pm + 2d \log\left(\frac{10}{\varepsilon}\right) + 44\bar{C}\sqrt{m} \log\left(\frac{epm}{44\bar{C}\sqrt{m}}\right) \\
&\quad \text{(Here } p = 8\sqrt{d}\varepsilon^{7/8} + \frac{3C_0}{\sqrt{m}} \leq \frac{11\bar{C}}{\sqrt{m}}) \\
&= -(8\sqrt{d}\varepsilon^{7/8} + \frac{3C_0}{\sqrt{m}})m + 2d \log(10m^{4/7}d^{4/7}) + 44\bar{C}\sqrt{m} \log\left(\frac{11e\bar{C}\sqrt{m}}{44\bar{C}\sqrt{m}}\right) \\
&\leq -8\sqrt{m} + 2.3d \log(md) \quad (m \geq 58) \\
&\leq -8\sqrt{m} + 7\sqrt{m} \quad (m \geq d^{2.4} \implies \sqrt{m} \geq 0.33d \log(md)) \\
&\leq -\sqrt{m} \\
&\leq \log(\gamma)
\end{aligned}$$

The last line holds for $m \geq d^{2.4}$ and $d \geq \log(4m/\gamma)^2 \geq \log(4/\gamma)^2$.

By the concentration of Lipschitz functions of Gaussians (Theorem B.3) and a union bound, we have

$$P\left(\exists S \in [m], |S| \leq 44\bar{C}\sqrt{m}, \left\| \frac{1}{\sqrt{m}} \sum_{s \in S} a_s \mathbf{w}_{s,0} \right\| \geq \sqrt{\frac{|S|}{m}}(\sqrt{d} + t)\right) \leq m^{44\bar{C}\sqrt{m}} e^{-\frac{t^2}{2}}$$

By setting $t = \sqrt{88\bar{C}\sqrt{m} \log(m) + 2 \log(\frac{8}{\gamma})}$, we get that with probability at least $1 - \gamma/8$,

$$\forall S \subset [m], |S| \leq 44\bar{C}\sqrt{m} : \left\| \frac{1}{\sqrt{m}} \sum_{s \in S} a_s \mathbf{w}_{s,0} \right\| \leq \sqrt{\frac{44\bar{C}}{\sqrt{m}}} \sqrt{88\bar{C}\sqrt{m} \log(m) + 2 \log(\frac{8}{\gamma})} + \sqrt{\frac{44\bar{C}d}{\sqrt{m}}}$$

$$\begin{aligned}
\forall S \subset [m], |S| \leq 44\bar{C}\sqrt{m} : & \left\| \sup_{\mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}} \frac{1}{\sqrt{m}} \sum_{i \in S} a_s \mathbf{w}_s \right\| \\
& \leq \left\| \frac{1}{\sqrt{m}} \sum_{i \in S} a_s \mathbf{w}_{s,0} \right\| + \left\| \sup_{\mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}} \frac{1}{\sqrt{m}} \sum_{i \in S} a_s \mathbf{v}_s \right\| \\
& \leq \sqrt{\frac{44\bar{C}}{\sqrt{m}}} \sqrt{88\bar{C}\sqrt{m} \log(m) + 2 \log(\frac{8}{\gamma})} + \sqrt{\frac{44\bar{C}d}{\sqrt{m}}} + \frac{44\bar{C}C_0}{\sqrt{m}} \\
& \leq 113\bar{C}C_0 \sqrt{\log(m) + \frac{2}{\sqrt{m}} \log(\frac{8}{\gamma})}
\end{aligned}$$

Note that for all $(r, \delta) \in N$, $\|\delta - \delta'\| \leq \varepsilon$,

$$\begin{aligned}
\sup_{\mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}} \mathbf{1}\{\text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top (\mathbf{x} + \delta)) < \text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top (\mathbf{x} + \delta'))\} & \leq X_s(\delta) \\
\sup_{\mathbf{v}_s \in \mathbb{R}^d, \|\mathbf{v}_s\| \leq \frac{C_0}{\sqrt{m}}} \mathbf{1}\{\text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top (\mathbf{x} + \delta)) > \text{sign}((\mathbf{w}_{s,0} + \mathbf{v}_s)^\top (\mathbf{x} + \delta'))\} & \leq X_s(\delta)
\end{aligned}$$

With probability at least $1 - \gamma$, we have for all δ, δ', r, r' with $\|\delta - \delta'\| \leq m^{-4/7}d^{-4/7}$ and $\|r - r'\| \leq m^{-4/7}d^{-4/7}$,

$$\begin{aligned}
|\Phi(r, \delta) - \Phi(r', \delta)| & \leq m^{-4/7}d^{-4/7} \left(\sqrt{md + 4m\sqrt{d \log(m/\gamma)}} + C_0 \right) \\
|\Phi(r, \delta) - \Phi(r, \delta')| & \leq 226\bar{C}C_0 \sqrt{\log(m) + \frac{2}{\sqrt{m}} \log(\frac{8}{\gamma})}
\end{aligned}$$

Combining this with (10) we obtain with probability at least $1 - \gamma$,

$$\begin{aligned}
& \sup_{(r, \delta) \in \Omega} \Phi(r, \delta) \\
& \leq 2\sqrt{2d \log(10m^{4/7}d^{4/7}) + \log(4/\gamma)} \left((4R\sqrt{\log(d)})^{1/4} + \sqrt{\frac{2d \log(10m^{4/7}d^{4/7}) + \log(4/\gamma)}{m}} \right) \\
& \quad + C_0 + 3\sqrt{\frac{1}{\sqrt{m}}(C_0 + \sqrt{\frac{\log(4/\gamma)}{2}})} \sqrt{d + 4\sqrt{d \log(8m/\gamma)}} \\
& \quad + 5(C_0 + \sqrt{\log(4/\gamma)}) \cdot \left(\sqrt{\log(4m/\gamma)} + 3 \right) + m^{-4/7}d^{-4/7} \left(\sqrt{md + 4m\sqrt{d \log(m/\gamma)}} + C_0 \right) \\
& \quad + 226\bar{C}C_0 \sqrt{\log(m) + \frac{2}{\sqrt{m}} \log(\frac{8}{\gamma})} \\
& \leq 6\sqrt{d \log(md)} \left((4R\sqrt{\log(d)})^{1/4} + \sqrt{\frac{6d \log(md)}{m}} \right) + 50d^{0.25} + 227\bar{C}C_0 \sqrt{\log(m) + \frac{2}{\sqrt{m}} \log(\frac{8}{\gamma})} \\
& \leq 9 \left(C_1 d^2 \log^2(md) \sqrt{\frac{\log(d)}{d}} \right)^{1/4} + \frac{15d \log(md)}{\sqrt{m}} + 327\bar{C}C_0 d^{0.25} \quad (R \leq \frac{C_1}{\sqrt{d}}) \\
& = o(\sqrt{d}) \quad (\text{holds when } \log^2(md) \sqrt{\frac{\log(d)}{d}} = o(1), \frac{d \log(md)}{\sqrt{m}} = o(\sqrt{d}) \implies m \leq O(\exp(d^{0.24})))
\end{aligned}$$

□

Proof of Corollary 3.2. From Theorem 3.1, we have that for any point $x \in \mathbb{S}^{d-1}$, with probability at least $1 - \gamma$, there exists an adversarial example for the neural network, with parameters (W, a) , at x for perturbation size R . Thus, with probability at least $1 - \gamma$, the robust error of W at x with perturbation R is one: $\ell_R(W; x, y) = 1$. This gives us that $\mathbb{E}_W[\inf_{W \in \mathcal{B}_{2,\infty}(W_0, \frac{C_0}{\sqrt{d}})} \ell_R(W; x, y)] \geq 1 - \gamma$. From

Markov's inequality,

$$\begin{aligned}
\mathbb{P} \left[\inf_{W \in \mathcal{B}_{2,\infty}(W_0, \frac{C_0}{\sqrt{d}})} L_R(W) < 0.9 \right] &= \mathbb{P} \left[1 - \inf_{W \in \mathcal{B}_{2,\infty}(W_0, \frac{C_0}{\sqrt{d}})} L_R(W) > 0.1 \right] \\
&\leq 10 \mathbb{E}_W \left[1 - \inf_{W \in \mathcal{B}_{2,\infty}(W_0, \frac{C_0}{\sqrt{d}})} L_R(W) \right] \\
&= 10 \mathbb{E}_x \mathbb{E}_W \left[1 - \inf_{W \in \mathcal{B}_{2,\infty}(W_0, \frac{C_0}{\sqrt{d}})} \ell_R(W; x, y) \right] \\
&\leq 10 \sup_x \mathbb{E}_W \left[1 - \inf_{W \in \mathcal{B}_{2,\infty}(W_0, \frac{C_0}{\sqrt{d}})} \ell_R(W; x, y) \right] \\
&\leq 10\gamma
\end{aligned}$$

where in the second equality we swap the expectations using Fubini's theorem. \square