# Adversarial Parameter Attack on Deep Neural Networks[*]

Lijia Yu, Yihan Wang, and Xiao-Shan Gao

Academy of Mathematics and Systems Science, Chinese Academy of Sciences,
Beijing 100190, China
University of Chinese Academy of Sciences, Beijing 100049, China

March 22, 2022

## Abstract

In this paper, a new parameter perturbation attack on DNNs, called adversarial parameter attack, is proposed, in which small perturbations to the parameters of the DNN are made such that the accuracy of the attacked DNN does not decrease much, but its robustness becomes much lower. The adversarial parameter attack is stronger than previous parameter perturbation attacks in that the attack is more difficult to be recognized by users and the attacked DNN gives a wrong label for any modified sample input with high probability. The existence of adversarial parameters is proved. For a DNN $\mathcal{F}_\Theta$ with the parameter set $\Theta$ satisfying certain conditions, it is shown that if the depth of the DNN is sufficiently large, then there exists an adversarial parameter set $\Theta_a$ for $\Theta$ such that the accuracy of $\mathcal{F}_{\Theta_a}$ is equal to that of $\mathcal{F}_\Theta$, but the robustness measure of $\mathcal{F}_{\Theta_a}$ is smaller than any given bound. An effective training algorithm is given to compute adversarial parameters and numerical experiments are used to demonstrate that the algorithms are effective to produce high quality adversarial parameters.

**Keyword**. Adversarial parameter attack, adversarial samples, robustness measurement, adversarial accuracy, mathematical theory for safe DNN.

## 1 Introduction

The deep neural network (DNN) [15] has become the most powerful machine learning method, which has been successfully applied in computer vision, natural language processing, and many other fields. Safety is a key desired feature of DNNs, which was studied extensively [1, 5, 42].

The most widely studied safety issue for DNNs is the adversarial sample attack [31], that is, it is possible to intentionally make small modifications to a sample, which are essentially imperceptible to the human eye, but the DNN outputs a wrong label or even any label given by the adversary. Existence of adversary samples makes the DNN vulnerable in safety-critical applications and many effective methods were proposed to develop more robust DNNs against adversarial attacks [20, 1, 5, 42]. However, it was shown that adversaries samples are inevitable for current DNN models in certain sense [6, 3, 27].

More recently, the parameter perturbation attacks [18, 43, 7, 33, 30, 37, 34, 35] were studied and shown to be another serious safety treat to DNNs. It was shown that by making small parameter

---

perturbations, the attacked DNN can give wrong or desired labels to specified input samples and still give the correct labels to other samples [18, 43, 34, 35].

In this paper, the adversarial parameter attack is proposed, in which small perturbations to the parameters of a DNN are made such that the attack to the DNN is essentially imperceptible to the user, but the robustness of the DNN becomes much lower. The adversarial parameter attack is stronger than previous parameter perturbation attacks in that not only the accuracy but also the robustness of DNNs are considered.

## 1.1 Contributions

Let $\mathcal{F}_\Theta$ be a DNN with $\Theta$ as the parameter set. A parameter perturbation $\Theta_a$ is called a set of *adversarial parameters* of $F_\Theta$ or $\Theta$, if the following conditions are satisfied 1) $\Theta_a$ is a small modification of $\Theta$, for instance $||\Theta_a - \Theta||_\infty \le \epsilon$ for a small positive number $\epsilon$; 2) the accuracy of $\mathcal{F}_{\Theta_a}$ over a distribution of samples is almost the same as that of $\mathcal{F}_\Theta$; 3) $\mathcal{F}_{\Theta_a}$ is much less robust than $\mathcal{F}_\Theta$, that is, $\mathcal{F}_{\Theta_a}$ has much more adversarial samples than $\mathcal{F}_\Theta$. It is clear that conditions 1) and 2) are to make the attack difficult to be recognized by the users and condition 3) is to make the new DNN less safe. The DNN obtained by the above attack is called an *adversarial DNN*, which has high accuracy but low robustness.

The existence of adversarial parameters is proved under certain assumptions. It is shown that if the depth of a trained DNN $\mathcal{F}_\Theta$ is sufficiently large, then there exist adversarial parameters $\Theta_a$ such that the accuracy of $\mathcal{F}_{\Theta_a}$ is equal to that of $\mathcal{F}_\Theta$, but the robustness measure of $\mathcal{F}_{\Theta_a}$ is as small as possible (refer to Corollaries 4.4 and 4.5). Since $\mathcal{F}_\Theta$ is a continuous function in $\Theta$, if $\Theta_a$ is an adversarial parameter for $\Theta$ then there exists a small sphere $S_a$ with $\Theta_a$ as center such that all parameters in $S_a$ are also adversarial parameters for $\Theta$. These results imply that adversarial parameters are inevitable in certain sense, similar to adversarial samples [6, 3, 27].

The existence of adversarial samples is usually demonstrated with numerical experiments, besides a few cases to be mention in the next section. As an application of adversarial parameters, we can construct DNNs which are guaranteed to have adversarial samples. For a trained DNN $\mathcal{F}_\Theta$ satisfying certain conditions, it is shown that there exist adversarial parameters $\Theta_a$ such that the accuracy of $\mathcal{F}_{\Theta_a}$ is equal to that of $\mathcal{F}_\Theta$, but $\mathcal{F}_{\Theta_a}$ has adversarial samples near a given normal sample (refer to Theorem 4.1), or the probability for $\mathcal{F}_{\Theta_a}$ to have adversarial samples over a distribution of samples is at least $1/2$ (refer to Theorem 4.2).

Finally, an effective training algorithm is given to compute adversarial parameters and numerical experiments are used to demonstrate that the algorithms are effective to produce high quality adversarial parameters for the networks VGG19 and Resnet56 on the CIFAR-10 dataset.

## 1.2 Related work

There exist vast literatures on adversarial attacks, which can be found in the survey papers [1, 5, 42]. We will focus on those which are closely related to our work.

**Parameter perturbation attacks.** Parameter perturbation attacks were given under different names such as fault injection attack, fault sneaking attack, stealth attack, and weight corruption. The fault injection attack [18] was first proposed by Liu et al, where it was shown that parameter perturbations can be used to misclassify one given input sample. In [7], it was shown that laser injection techniques can be used as a successful fault injection attack in real-world applications. In [43], the fault sneaking attack was proposed, where multiple input samples were misclassified and

2

other samples were still given the correct label. In [37], lower bounds were given for parameter perturbations under which the network still gives the correct label for a given sample. In [33], upper bounds were given for the changes of the pairwise class margin function and the Rademacher complexity against parameter perturbations and new loss functions were given to obtain more robust networks. In [30], the maximum change of the loss function over given samples was used as an indicator to measure the robustness of DNNs against parameter perturbations and gradient decent methods were used to compute the indicator. In [34, 35], the stealth attack which can guaranteed to make the attacked DNN gives a desired label for a sample outside of the validation set and keep correct labels for samples in the validation set. The stealth attack has the form $\mathcal{F} + \mathcal{U}$, where $\mathcal{F}$ is the DNN to be attacked and $\mathcal{U}$ is a DNN with one hidden layer.

The adversarial parameter attack proposed in this paper is stronger than previous parameter perturbation attacks by in the following aspects. First, by keeping the accuracy and eliminating the robustness, the adversarial parameter attack is more difficult to be recognized, because the attached DNN performs almost the same as the original DNN on the test set. Second, by reducing the robustness of the DNN, the attacked DNN gives a wrong label for any modified input sample with high probability, while previous parameter attacks usually misclassify certain given samples. Finally, we prove the existence of adversarial parameters under reasonable assumptions.

**Mathematical theories of adversarial samples.** Existence of adversarial samples were usually demonstrated with numerical experiments, and mathematical theories were desired. In [6], it was proved that for DNNs with a fixed architecture, there exist uncountable classification functions and distributions of samples such that adversarial samples always exist for any successfully trained DNN with the given architecture and the sample distribution. In the stealth attack [34, 35], it was proved that there exist attached DNNs which give a desired label for a sample outside of the validation set by modifying the DNN. In this paper, we show that by making small perturbations to the parameters of the DNN, the DNN has adversarial samples with high probability.

Theories for certified robustness of DNNs were given in several aspects. Let $x$ be a sample such that the DNN $\mathcal{F}$ gives the correct label. Due to the continuity of the DNN function, a sphere with $x$ as center does not contain adversarial samples if its radius is sufficiently small, which is called a *robust sphere* of $x$. In [12], lower bounds for the robustness radius were computed and used to enhance the robustness of the DNN. In [24], for shallow networks, the upper bounds of the changes of the network under sample input perturbations were given and use to obtain more robust DNNs. In [8], the random smoothing method was proposed and lower bounds for the radius of the robust spheres was given. In [39], lower bounds for the average radius of robust spheres for a distribution of samples are given. Universal lower bounds on stability in terms of the dimension of the domain of the classification function were also given in [27, 34]. However, these bounds are usually inverse-exponentially dependent on the depth of the DNN, which are very small for deep networks in real world applications. In [40], the information-theoretically safe bias classifier was introduced by making the gradient of the DNN random.

**Algorithms to train robust DNNs.** Many effective methods were proposed to train more robust DNNs to defend adversarial samples [1, 5, 42, 38]. Methods to train DNNs which are more robust against parameter perturbation attacks were also proposed [18, 43]. The adversarial training method proposed by Madry et al [20] can reduce adversarial samples significantly, where the value of the loss function of the worst adversary in a small neighborhood of the training sample is minimized. In this paper, the idea of adversarial training is used to compute adversarial parameters.

## 2 Adversarial parameters

In this section, we define the adversarial parameters and give a measurement for the quality of the adversarial parameters.

### 2.1 Adversarial parameters of DNNs

Let us consider a standard DNN. Denote $\mathbb{I} = [0,1] \subset \mathbb{R}$ and $[n] = \{1, \ldots, n\}$ for $n \in \mathbb{N}_+$. In this paper, we assume that $\mathcal{F} : \mathbb{I}^n \to \mathbb{R}^m$ is a classification DNN for $m$ objects, which has $L$ hidden-layers, all hidden-layers use Relu as the activity function, and the output layer does not have activity functions. $\mathcal{F}$ can be written as

$$
\begin{aligned}
&x_0 \in \mathbb{I}^n, n_0 = n, n_{L+1} = m; \\
&x_l = \text{Relu}(W_l x_{l-1} + b_l) \in \mathbb{R}^{n_l}, W_l \in \mathbb{R}^{n_l \times n_{l-1}}, b_l \in \mathbb{R}^{n_l}, l \in [L]; \\
&\mathcal{F}(x_0) = x_{L+1} = W_{L+1} x_L + b_{L+1}, W_{L+1} \in \mathbb{R}^{m \times n_L}, b_{L+1} \in \mathbb{R}^m.
\end{aligned}
\tag{1}
$$

Denote $\Theta = \{W_l, b_l\}_{l=1}^{L+1} \in \mathbb{R}^k$ to be the parameter set of $\mathcal{F}$ and $\mathcal{F}$ is denoted as $\mathcal{F}_\Theta$ if the parameters need to be mentioned explicitly, where $k = \sum_{l=1}^{L+1} n_l(n_{l-1} + 1)$.

Let $\mathcal{F}_\Theta$ be a trained network with the parameter set $\Theta$. Then a new parameter set $\Theta_a$ is called a set of *adversarial parameters* of $\Theta$ if 1) $\Theta_a$ is a small perturbation of $\Theta$; 2) the accuracy of $\mathcal{F}_{\Theta_a}$ is almost the same as that of $\mathcal{F}_\Theta$; 3) $\mathcal{F}_{\Theta_a}$ is much less robust comparing to $\mathcal{F}_\Theta$, that is, $\mathcal{F}_{\Theta_a}$ has more adversarial samples than $\mathcal{F}_\Theta$.

We assume that the objects to be classified satisfy a distribution $D_x$ in $\mathbb{R}^n$, and a sample $x \sim D_x$ is called a *normal sample*. Let $\Theta \in \mathbb{R}^k$ be the parameter set of a trained network $\mathcal{F}_\Theta : \mathbb{I}^n \to \mathbb{R}^m$. For $x \sim D_x$, denote $l_x$ to be the label of $x$ and $\widehat{\mathcal{F}}_\Theta$ to be the classification result of $\mathcal{F}_\Theta$. Then the accuracy of $\mathcal{F}_\Theta$ for the normal samples is

$$
A(\mathcal{F}_\Theta, D_x) = P_{x \sim D_x}(\widehat{\mathcal{F}}_\Theta(x) = l_x).
\tag{2}
$$

In order to measure the quality of adversarial parameters, we need a robustness measure $R(\mathcal{F}_\Theta, D_x)$ of $\mathcal{F}_\Theta$ for the normal samples. There exist several definitions for $R(\mathcal{F}_\Theta, D_x)$ [20, 39]. In this paper, two kinds of robustness measures are used.

We first give two robust measures of $\mathcal{F}_\Theta$ on a given sample $x_0$. The *robustness radius* of $x_0$ under the $L_p$ norm for $p \in \mathbb{R}_+ \cup \{\infty\}$ is defined to be

$$
R_1(\mathcal{F}_\Theta, x_0) = \max\{\zeta \in \mathbb{R}_+ \mid \widehat{\mathcal{F}}(x) = l_x, \forall x \text{ s.t. } ||x - x_0||_p \leq \zeta\}.
\tag{3}
$$

If $\widehat{\mathcal{F}}_\Theta(x_0) \neq l_{x_0}$, then the robustness radius of $x_0$ is zero. It is difficult to have good estimation to the robustness radius, and the following approximation to the robust radius under $L_p$ norm [12] is often used

$$
R_2(\mathcal{F}_\Theta, x_0) = \min_{l \in [m], l \neq l_x} \left\{ \frac{|\mathcal{F}_{l_x}(x_0) - \mathcal{F}_l(x_0)|}{||\nabla \mathcal{F}_{l_x}(x_0) - \nabla \mathcal{F}_l(x_0)||_q} I(\mathcal{F}_{l_x}(x_0) > \mathcal{F}_l(x_0)) \right\}
\tag{4}
$$

where $\mathcal{F}_l(x_0)$ is the $l$-th coordinate of $\mathcal{F}(x_0)$, $\nabla(\mathcal{F}_l(x)) = \frac{\nabla \mathcal{F}_l(t)}{\nabla t}|_{t=x}$, $p, q \in \mathbb{R}_+\{\infty\}$ satisfy $1/q + 1/p = 1$ ($p = 0(\infty)$ iff $q = \infty(0)$), and $I(t) = 1$ if $t$ it true or $I(t) = 0$ otherwise.

For a distribution $D_x$ of samples, we define two global robust measures corresponding to $R_1$ and $R_2$. The adversarial accuracy can be used as $R(\mathcal{F}_\Theta, D_x)$. For $\epsilon \in \mathbb{R}_+$ and $p \in \mathbb{R}_+ \cup \{\infty\}$, the

adversarial accuracy of $\mathcal{F}_\Theta$ is

$$
\begin{aligned}
R_3(\mathcal{F}_\Theta, D_x, \epsilon) &= P_{x \sim D_x}(\epsilon \le R_1(\mathcal{F}_\Theta, x)) \\
&= P_{x \sim D_x}(\widehat{\mathcal{F}}_\Theta(x') = l_x, \forall x' \text{ s.t. } ||x' - x||_p \le \epsilon).
\end{aligned}
\tag{5}
$$

Corresponding to $R_2$ in (4), we have the following global robustness measure

$$
R_4(\mathcal{F}, D_x) = \int_{x \sim D_x} R_2(\mathcal{F}, x) \mathrm{d}x.
\tag{6}
$$

We now define a measurement for an adversarial parameter set using the accuracy and robustness of $\mathcal{F}$.

**Definition 2.1.** *Let $\Theta_a$ be an adversarial parameter set of $\Theta$, $R(\mathcal{F}, D_x)$ a robustness measure of $\mathcal{F}$ for normal samples, and*

$$
\begin{aligned}
P_{x \sim D_x}(\mathcal{F}_{\Theta_a}(x) = l_x) &= \gamma_1 P_{x \sim D_x}(\mathcal{F}_\Theta(x) = l_x) \\
R(\mathcal{F}_{\Theta_a}, D_x) &= \gamma_2 R(\mathcal{F}_\Theta, D_x).
\end{aligned}
\tag{7}
$$

*Then the* adversarial rate *of $\Theta_a$ is defined to be $\overline{\gamma_1}(1 - \overline{\gamma_2})$, where $\overline{\gamma} = \min\{\gamma, 1\}$.*

In general, we have $\gamma_1 \le 1$ and $\gamma_2 \le 1$. The value of $\gamma_1$ measures the ability of $\Theta_a$ to keep the accuracy of $\mathcal{F}_\Theta$ on normal samples, and if $\gamma_1$ is large then the attack is more difficult to be detected. The value of $1 - \gamma_2$ measures the ability of $\Theta_a$ to break the robustness of $\mathcal{F}_\Theta$, and if $1 - \gamma_2$ is large then the parameter attack is more powerful. Hence, the adversarial rate $\overline{\gamma_1}(1 - \overline{\gamma_2})$ measures the quality of the adversarial parameter attack in that if the adversarial rate is larger then the adversarial parameter attack is better. If $\overline{\gamma_1}(1 - \overline{\gamma_2})$ achieves its maximal value 1, then $\gamma_1 = 1$ and $\gamma_2 = 0$ and the adversarial parameter attack is a perfect attack in that the attack does not change the accuracy of $\mathcal{F}$, but totally destroys the robustness of $\mathcal{F}$.

**Remark 2.1.** *In order to make the attack very hard to be detected, we can give a lower bound $\gamma_{\text{low}}$ to $\gamma_1$. If $\gamma_1 < \gamma_{\text{low}}$, we consider $\Theta_a$ to be a failed attack.*

## 2.2   Adversarial parameter attacks for other purposes

According to the requirements of specific applications, we may define other types of adversarial parameter attacks.

The adversarial parameters defined in section 2.1 are for all the samples. In certain applications, it is desired to make the network less robust on one specific class of samples, which motivates the following definitions.

The simplest case is adversarial parameters for a given sample. A small perturbation $\Theta_a$ of $\Theta$ is called *adversarial parameters* for a given sample $x_0$, if $\mathcal{F}_{\Theta_a}$ gives the correct label to $x_0$ and has adversarial samples of $x_0$ in $S_\infty(x_0, \epsilon) = \{x \,|\, |x - x_0|_\infty \le \epsilon\}$ for a given $\epsilon \in \mathbb{R}_+$. Let $R(\mathcal{F}_\Theta, x_0)$ be a measure of robustness of $\mathcal{F}_\Theta$ at sample $x_0$, and

$$
R(\mathcal{F}_{\Theta_a}, x_0) = \gamma R(\mathcal{F}_\Theta, x_0).
\tag{8}
$$

Then the adversarial rate of $\Theta_a$ is defined to be $1 - \overline{\gamma}$.

We can also break the stability for samples with a given label. A small perturbation $\Theta_a$ of $\Theta$ is called *adversarial parameters* for samples with a given label $l_0$, if $\mathcal{F}_{\Theta_a}$ keeps the accuracy for all

normal samples and the robustness for normal samples whose label is not $l_0$, but break the robustness of samples with label $l_0$. Let $\alpha = P_{x \sim D_x}(\widehat{\mathcal{F}}_\Theta(x) = l_x)$ and $\beta = P_{x \sim D_x}(\widehat{\mathcal{F}}_\Theta(x') = l_x, \forall x' \in S_p(x, \epsilon))$. For such adversarial parameters $\Theta_a$, let

$$P_{x \sim D_x}(\widehat{\mathcal{F}}_{\Theta_a}(x) = l_x) = \gamma_1 \alpha \tag{9}$$
$$P_{x \sim D_x}(\widehat{\mathcal{F}}_{\Theta_a}(x') = l_x, \forall x' \in B_p(x, \epsilon) \mid l_x \neq y_0) = \gamma_2 \beta$$
$$P_{x \sim D_x}(\widehat{\mathcal{F}}_{\Theta_a}(x') = l_x, \forall x' \in B_p(x, \epsilon) \mid l_x = y_0) = \gamma_3 \beta.$$

Then the adversarial rate of $\Theta_a$ is defined as $\overline{\gamma}_1 \overline{\gamma}_2 (1 - \overline{\gamma}_3)$.

Finally, instead of breaking the robustness of samples with label $y_0$, we can break the accuracies for samples with label $y_0$. Such adversarial parameters are called *direct adversarial parameters*. Let $\Theta_a$ be a direct adversarial parameter set and

$$P_{x \sim D_x}(\widehat{\mathcal{F}}_{\Theta_a}(x) = l_x \mid l_x \neq y_0) = \gamma_1 \alpha \tag{10}$$
$$P_{x \sim D_x}(\widehat{\mathcal{F}}_{\Theta_a}(x') = l_x, \forall x' \in B_p(x, \epsilon) \mid l_x \neq y_0) = \gamma_2 \beta$$
$$P_{x \sim D_x}(\widehat{\mathcal{F}}_{\Theta_a}(x) = l_x \mid l_x = y_0) = \gamma_3 \alpha.$$

Then the adversarial rate is defined as $\overline{\gamma}_1 \overline{\gamma}_2 (1 - \overline{\gamma}_3)$. The above definition is similar to the attacks in [18, 43], but robustness is considered as an extra objective.

# 3   Algorithm

In this section, we give algorithms to compute adversarial parameters.

## 3.1   Compute adversarial parameters under $L_\infty$ norm

We formulate the adversarial parameter attack for a trained DNN $\mathcal{F}_\Theta$ under the $L_p$ norm as the following optimization problem for a given $\zeta \in \mathbb{R}_+$.

$$\max_{\Theta_a \in \mathbb{R}^k, ||\Theta_a - \Theta||_p \leq \zeta} A(\mathcal{F}_{\Theta_a}, D_x) / R(\mathcal{F}_{\Theta_a}, D_x) \tag{11}$$

where $A(\mathcal{F}_{\Theta_a}, D_x)$ and $R(\mathcal{F}_{\Theta_a}, D_x)$ are the accuracy and a robustness measure for $\mathcal{F}_{\Theta_a}$ over a distribution sample $D_x$.

**Remark 3.1.** *Theoretically, the adversarial parameter attack should be a multi-objective optimization problem, that is to maximize the accuracy and to minimize the robustness. But, such an optimization problem is difficult to solve.*

**Remark 3.2.** *According to (11), the adversarial rate seems better to be defined as $\gamma_1/\gamma_2$, which is possible but not as good as the one in Definition 2.1 for the following reasons. The adversarial rate $\overline{\gamma_1}(1 - \overline{\gamma_2})$ has the optimal value 1 and gives a more intuitive view to see the quality of the adversarial parameters.*

In the rest of this section, we show how to change formula (11) to an effective algorithm to compute $L_\infty$ norm adversarial parameters using the robustness measure in (5). We first show how to compute the robustness in (5) explicitly. We use the adversarial training [20] to do that, which

is the most effective way to find adversarial samples. For a sample $x$ and a small number $\varepsilon \in \mathbb{R}_+$, we first compute

$$\chi_0 = \arg\max_{\chi \in \mathbb{R}^k, ||\chi||_0 < \varepsilon} L_{\mathrm{CE}}(\mathcal{F}_\Theta(x + \chi), l_x)$$

with PGD [20] and then use

$$L_{\mathrm{AT}}(x, \Theta) = L_{\mathrm{CE}}(\mathcal{F}_\Theta(x + \chi_0), l_x) \tag{12}$$

to measure the robustness of $\mathcal{F}_\Theta$ at $x$.

We need a training set $T$ to find the adversarial parameters. The training procedure consists of two phases. In the first pre-training phase, the loss function

$$-\sum_{x \in T} L_{\mathrm{AT}}(x, \Theta) \tag{13}$$

is used to reduce the adversarial accuracy of $\mathcal{F}_\Theta$. In the second main training phase, the loss function

$$\frac{\sum_{x \in T} L_{\mathrm{CE}}(\mathcal{F}_\Theta(x), l_x)}{\sum_{x \in T} L_{\mathrm{AT}}(x, \Theta)} \tag{14}$$

is used to promote the accuracy and keep the low-level of adversarial accuracy of $\mathcal{F}_\Theta$, which corresponds to formula (11).

We will compute a more general $L_\infty$ norm parameter perturbation. Let $\Delta \in \mathbb{R}_+^k$ and $\Delta_i$ the $i$-th coordinate of $\Delta$. Then the $L_\infty$ parameter perturbation will be found in

$$B_\infty(\Theta, \Delta) = \{\Theta_a \in \mathbb{R}^k \,|\, |\Theta_a - \Theta|_i \leq \Delta_i, \ \forall \ i \in [k]\}.$$

It is clear that the usual $L_\infty$ norm parameter perturbation is a special case of the above general case. We use this general form, because we want to include more types of parameter perturbations which are given in section 5.2. A sketch of the algorithm is given below.

---

**Algorithm 1** Attack under $L_\infty$ norm

---

**Require:**
  The parameter set $\Theta$ of $\mathcal{F}$;
  The hyper-parameters: $\alpha \in \mathbb{R}_+$, $\Delta \in \mathbb{R}_+^k$, $n_1, n_2 \in \mathbb{N}$;
  A training set $T$.
**Ensure:** An adversarial parameter set $\Theta_a$ in $B_\infty(\Theta, \Delta)$.
  Let $i = 0$, $\Theta_a = \Theta$.
  For all $i \in [n_1 + n_2]$:
    If $i < n_1$:
      $L = -\sum_{x \in T} \frac{1}{|T|} L_{\mathrm{AT}}(x, \Theta_a)$.
    Else:
      $L = \frac{\sum_{x \in T} L_{\mathrm{CE}}(\mathcal{F}_{\Theta_a}(x), l_x)}{\sum_{x \in T} L_{\mathrm{AT}}(x, \Theta_a)}$.
    $\widetilde{\Theta} = \Theta_a + \alpha \bigtriangledown L$.
    $\Theta_a = \mathrm{Proj}(\widetilde{\Theta}, B_\infty(\Theta, \Delta))$.
  Output: $\Theta_a$.

---

**Remark 3.3.** *We give more details for the algorithm.*

*(1).* $\mathrm{Proj}(\widetilde{\Theta}, B_\infty(\Theta, \Delta))$ *maps* $\widetilde{\Theta}$ *into* $B_\infty(\Theta, \Delta)$ *as follows: for* $i \in [k]$:

7

If $\widetilde{\Theta}_i > \Theta_i + \Delta_i$, $\text{Proj}(\widetilde{\Theta}, B_\infty(\Theta, \Delta))_i = \Theta_i + \Delta_i$;

If $\widetilde{\Theta}_i < \Theta_i - \Delta_i$, $\text{Proj}(\widetilde{\Theta}, B_\infty(\Theta, \Delta))_i = \Theta_i - \Delta_i$;

If $\widetilde{\Theta}_i + \Delta_i > \Theta_i > \widetilde{\Theta}_i - \Delta_i$, $\text{Proj}(\widetilde{\Theta}, B_\infty(\Theta, \Delta))_i = \Theta_i$.

(2). We will reduce the training steps $\alpha$ with the training going.

## 3.2 Algorithms for other kinds of adversarial parameters

The algorithm to find adversarial parameters under other norms and robustness measures can be developed similarly. In what below, we show how to compute adversarial parameters under $L_0$ norm, which is different from other cases. The overall algorithm is similar to Algorithm 1, except we use a new method to update the parameters. Suppose that $\Theta = \{W_l, b_l\}_{l=1}^{L}$ is the parameter to be updated and the value of $L$ in Algorithm 1 is found. We will show how to update the parameters. We only change some weight matrices $W_l$ as follows.

- Randomly select two entries $w_1$ and $w_2$ of $W_l$ until $(\frac{\nabla L}{\nabla w_1} - \frac{\nabla L}{\nabla w_2})(w_1 - w_2) > 0$ is satisfied.

- Exchange $w_1$ and $w_2$ in $W_l$ to obtain the new parameters.

It is clear that the change will make $L$ become smaller. In total, we update a given number of weight matrices, and for each such matrix, we change a given percentage of its entries. The details of the algorithm are omitted. Note that the above parameter perturbation keeps the sparsity and the values of the entries of the weight matrices. As a consequence the Proj operator in the algorithm can be taken as the identity map.

The adversarial parameters defined in section 2.2 can also be obtained similarly. For instance, to compute the adversarial parameters for one sample $x$, we just need to let $T$ in Algorithm 1 to be $T = \{x\}$.

To compute adversarial parameters for samples with a given label $l_0$, by (9) we can use the following loss function

$$\frac{\sum_{x \in T} L_{\text{CE}}(\mathcal{F}_{\Theta_a}(x), l_x) + \sum_{x \in T \ \& \ l_x \neq l_0} L_{\text{AT}}(x, \Theta_a)}{\sum_{x \in T \ \& \ l_x = l_0} L_{\text{AT}}(\mathcal{F}_{\Theta_a}(x), l_x)} \tag{15}$$

to increase the robustness and accuracy of samples whose labels are not $l_0$ and to reduce the accuracy for samples with labels $l_0$.

To compute direct adversarial parameters for samples with label $l_0$, by (10) we can use the following loss function

$$\frac{\sum_{x \in T \ \& \ l_x \neq l_0} (L_{\text{AT}}(x, \Theta_a) + L_{\text{CE}}(\mathcal{F}_{\Theta_a}(x), l_x))}{\sum_{x \in T \ \& \ l_x = l_0} L_{\text{CE}}(\mathcal{F}_{\Theta_a}(x), l_x)} \tag{16}$$

to increase the robustness and accuracy of samples whose labels are not $l_0$ and to reduce the accuracy for samples with labels $l_0$.

# 4 Existence of adversarial parameters

In this section, we will show that adversarial parameters with high adversarial rates exist under certain conditions.

## 4.1 Adversarial parameters to achieve low adversary accuracy

In this section, we use the robustness radius in (3) and the adversarial accuracy in (5) as the robust measures, and hence existence of adversarial parameters implies low adversary accuracies.

We introduce several notations. Let $||x||_{-\infty} = \min_{i \in [n]} \{|x_i|\}$ for $x \in \mathbb{R}^n$, and $||W||_{-\infty,2} = \min_{i \in [a]} \{||W^{(i)}||_2\}$ for $W \in \mathbb{R}^{a \times b}$, where $W^{(i)}$ is the $i$-th row of $W$. If $\mathcal{F}$ is a network, we use $\mathcal{F}_i(x)$ to denote the $i$-th coordinate of $\mathcal{F}(x)$.

In this section, we consider the following network $\mathcal{F}_\Theta : \mathbb{I}^n \to \mathbb{R}^m$ with one hidden layer

$$\mathcal{F}(x) = W_2 \text{Relu}(W_1 x + b_1) + b_2, \tag{17}$$

where $W_1 \in \mathbb{R}^{n_1 \times n}, b_1 \in \mathbb{R}^{n_1}, W_2 \in \mathbb{R}^{m \times n_1}, b_2 \in \mathbb{R}^m$. $\Theta = \{W_i, b_i\}_{i=1}^2 \in \mathbb{R}^k$ is the parameter set of $\mathcal{F}_\Theta$, where $k = (n + m + 1)n_1 + m$.

The network defined in (17) has just one hidden-layer. We will show that when the width of its hidden-layer is large enough, adversarial parameters exist under with certain conditions.

We will consider $L_\infty$ adversarial parameters. For $\gamma \in \mathbb{R}_+$, the hypothetical space for the adversarial networks of $\mathcal{F}_\Theta$ is

$$\mathcal{H}_\gamma(\Theta) = \{\mathcal{F}_{\Theta_a} \,|\, ||\Theta_a - \Theta||_\infty < \gamma\}.$$

The following theorem shows the existence of adversarial parameters for a given sample $x_0$. The proof of the theorem is given in section 6.1.

**Theorem 4.1.** *Let $\mathcal{F}_\Theta$ be a trained network with structure in (17), which gives the correct label $l_{x_0}$ for a sample $x_0$. Further assume the following conditions.*

$C_1$. *Let $a, A \in \mathbb{R}_+$ such that $|\mathcal{F}_i(x) - \mathcal{F}_j(x)| < A$ for all $i, j \in [m]$ and $x \in S_\infty(x_0, a) = \{x \,|\, ||x - x_0||_\infty \le a\}$.*

$C_2$. *$||W_2^{(i)} - W_2^{(j)}||_{-\infty} > c$ for all $i, j \in [m], i \ne j$.*

$C_3$. *At least $\eta n_1$ coordinates of $|\text{Relu}(W_1 x + b_1)|$ are bigger than $b$, where $\eta \in (0, 1)$ and $b \in \mathbb{R}_+$.*

*For $\gamma, \epsilon \in \mathbb{R}_+$ such that $\epsilon < a$, if $n_1 > \frac{2A}{\min\{\epsilon \gamma (n-1), b\} c \eta}$, then there exists an $\mathcal{F}_{\Theta_a} \in \mathcal{H}_\gamma(\Theta)$ such that $\widehat{\mathcal{F}}_{\Theta_a}(x_0) = l_{x_0}$ and $\mathcal{F}_{\Theta_a}$ has adversarial samples to $x_0$ in $S_\infty(x_0, \epsilon)$.*

We have the following observations from Theorem 4.1.

**Corollary 4.1.** *If the robustness radius in (3) is used as the robustness measure, then the adversarial rate of $\mathcal{F}_{\Theta_a}$ in Theorem 4.1 is bigger than $1 - \frac{\epsilon}{a}$.*

**Remark 4.1.** *From Theorem 4.1, if the width of $\mathcal{F}_\Theta$ is sufficiently large, then $\mathcal{F}_\Theta$ has adversarial parameters which is as close as possible to $\Theta$ and $\mathcal{F}_{\Theta_a}$ has adversarial samples which are as close as possible to $x_0$.*

**Remark 4.2.** *Since $\mathcal{F}_\Theta$ is a continuous function in $\Theta$, if $\Theta_a$ is an adversarial parameter set for $\Theta$ then there exists a small sphere $S_a$ with $\Theta_a$ as center such that all parameters in $S_a$ are also adversarial parameters for $\Theta$.*

From the above two remarks, we may say that adversarial parameters are inevitable in this case.

The following theorem shows that when $n_1$ is large enough, adversarial parameters exist for a distributions with high probability. The proof of the theorem is given in section 6.1.

**Theorem 4.2.** *Let $\mathcal{F}_\Theta$ be a trained DNN with structure in (17) and $S \subset \mathbb{I}^n$ the set of normal samples. Further assume the following conditions.*

$C_1$. *Let $A, a \in \mathbb{R}_+$ such that $|\mathcal{F}_i(x) - \mathcal{F}_j(x)| < A$ for all $i, j \in [m]$ and $x \in \bigcup_{x_0 \in S} S_\infty(x_0, a)$.*

$C_2$. $||W_2^{(i)} - W_2^{(j)}||_{-\infty} > c$ *for all $i, j \in [m]$, $i \neq j$, where $c \in \mathbb{R}_+$.*

$C_3$. *For all $x \in S$, at least $\eta n_1$ coordinates of $|\mathrm{Relu}(W_1 x + b_1)|$ are bigger than $b$, where $\eta \in (0,1)$ $b \in \mathbb{R}_+$.*

$C_4$. *The dimension of $S$ is lower than $n - m$.*

*For $\epsilon, \gamma \in \mathbb{R}_+$ such that $\epsilon < a$, if $n_1 > \frac{2A}{\min\{\epsilon\gamma/m, b\}c\eta}$, then there exists an $\mathcal{F}_{\Theta_a} \in \mathcal{H}_\gamma(\Theta)$ such that the accuracy of $\mathcal{F}_{\Theta_a}$ over $D_x$ is greater than or equal to that of $\mathcal{F}_\Theta$ and*

$$P_{x \sim D_x}(\mathcal{F}_{\Theta_a} \text{ has an adversarial sample of } x \text{ in } S_\infty(x, \epsilon)) \geq 0.5.$$

**Corollary 4.2.** *Let the adversary accuracy of $\mathcal{F}_\Theta$ be $\theta = R_3(\mathcal{F}_\Theta, D_x, \epsilon)$, then Theorem 4.2 implies that there exists adversarial parameters $\Theta_a$ of $\Theta$ with adversarial rate at least $1 - 0.5/\theta$.*

**Remark 4.3.** *The conditions of Theorems 4.1 and 4.2 can be satisfied for most DNNs. The parameters $A$ and $a$ in Condition $C_1$ are clearly exist. Since the training procedure usually terminates near a local minimum of the loss function, the weights can be considered as random values [39], and hence conditions $C_2$ and $C_3$ can be satisfied. For practical examples such as MNIST and CIFAR-10, condition $C_4$ is clearly satisfied.*

## 4.2 Adversarial parameters for DNNs

In this section, we consider networks of the following form

$$
\begin{aligned}
& x_0 \in \mathbb{I}^n, n_{L+1} = m; \\
& x_l = \mathrm{Relu}(W_l x_{l-1}) \in \mathbb{R}^n, W_l \in \mathbb{R}^{n \times n}, l \in [L]; \\
& \mathcal{F}(x_0) = x_{L+1} = W_{L+1} x_L \in \mathbb{R}^m, W_{L+1} \in \mathbb{R}^{m \times n}.
\end{aligned}
\tag{18}
$$

Let $\Theta = \{W_i\}_{i=1}^{L+1}$ be the parameters and $\Theta \in \mathbb{R}^k$, where $k = Ln^2 + mn$. We use $\mathcal{F}_\Theta^l(x)$ to represent the output of the $l$-th layer of $\mathcal{F}_\Theta(x)$ where $l \in [L]$. We will show that, when $L$ becomes big, adversarial parameters exist.

We first prove the existence of adversarial parameters for a given sample. We use the following robustness measure for network $\mathcal{F}$ at a sample $x$

$$R(\mathcal{F}, x) = \min_{l \neq l_x}\{\frac{|\mathcal{F}_{l_x}(x) - \mathcal{F}_l(x)|^2}{||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_l(x))||_2^2} I(\mathcal{F}_{l_x}(x) > \mathcal{F}_l(x))\}.$$

It is easy to see that this is the square of $R_2(\mathcal{F}, x)$ in (4) with $p = 2$.

**Theorem 4.3.** *Let $\mathcal{F}_\Theta$ be a trained network with structure in (18), which gives the correct label for a sample $x_0 \in \mathbb{R}^n$. Further assume the following conditions.*

$C_1$. $||\frac{\nabla \mathcal{F}_i(t)}{\nabla t}|_{t=x_0}||_2 < \sqrt{A}$ *for $i \in [m]$.*

$C_2$. $||\frac{\nabla \mathcal{F}^l(t)}{\nabla t}|_{t=x_0}||_{-\infty, 2} > b$ *for $l \in [L]$.*

10

$C_3$. $||\frac{\nabla \mathcal{F}_i(t) - \nabla \mathcal{F}_j(t)}{\nabla \mathcal{F}^l(t)}|_{t=x_0}||_{-\infty} > c$ for $i, j \in [m]$, $i \neq j$ and $l \in [L]$.

$C_4$. For $l \in [L]$, $\frac{\nabla \mathcal{F}^l(t)}{\nabla t}|_{t=x_0}$ has a column $L_l$ such that the angle between $L_l$ and $\mathcal{F}^l(x_0)$ is bigger than $\alpha$ and smaller than $\pi - \alpha$, where $\alpha \in [0, \pi/2]$.

Then for $\gamma \in \mathbb{R}_+$, there exists an $\mathcal{F}_{\Theta_a} \in \mathcal{H}_\gamma(\Theta)$ such that $\mathcal{F}_{\Theta_a}(x_0) = l_{x_0}$ and

$$R(\mathcal{F}_{\Theta_a}, x_0) \leq (1 - \eta)R(\mathcal{F}_\Theta, x_0)$$

where $\eta = \frac{\gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)}{4A + \gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)}$. In other words, there exists an $\mathcal{F}_{\Theta_a}$ with adversarial rate $\geq \eta$.

The proof of the theorem is given in section 6.2. As a consequence of Theorem 4.3, there exist adversarial parameters for sample $x_0$, whose robustness measure is as small as possible.

**Corollary 4.3.** For $\rho \in (0, 1)$, if $L \geq \frac{4(1-\rho)A}{\rho(\gamma \sin(r)cb)^2} + 1$, then the adversarial rate of $\mathcal{F}_{\Theta_a}$ is $\geq 1 - \rho$.

**Corollary 4.4.** For $\tau \in (0, 1)$ satisfying $\tau < R(\mathcal{F}, x_0)$, if $L > \frac{4A(R(\mathcal{F}, x_0)/\tau - 1)}{(\gamma \sin(r)cb)^2} + 1$, then $R(\mathcal{F}_{\Theta_a}, x_0) \leq \tau$.

To find adversarial parameters for samples under a distribution $D_x$, we use the following robustness measure for $\mathcal{F}$:

$$R(\mathcal{F}, D_x) = \frac{\int_{x \sim D_x} \min_{j \neq l_x}\{||\mathcal{F}_{l_x}(x) - \mathcal{F}_j(x)||_2^2 I(\mathcal{F}_{l_x}(x) > \mathcal{F}_j(x))\}\mathrm{d}x}{\int_{x \sim D_x} \max_{j \neq l_x}\{||\nabla \mathcal{F}_{l_x}(x) - \nabla \mathcal{F}_j(x)||_2^2\}\mathrm{d}x}.$$

This is a variant of $R_4(\mathcal{F}, D_x)$ in (6) with $p = 2$. The following theorem shows that adversarial parameters exist for this robustness measure. The proof of the theorem is given in section 6.2.

**Theorem 4.4.** Let $\mathcal{F}_\Theta$ be a trained DNN with structure in (17) and $S \subset \mathbb{I}^n$ the set of normal samples satisfying distribution $D_x$. Further assume the following conditions.

$C_1$. $||\frac{\nabla \mathcal{F}_i(t)}{\nabla t}|_{t=x}||_2 < \sqrt{A}$ for all samples $x \in S$ and $i \in [m]$.

$C_2$. $P_{x \sim D_x}(\forall l \neq l_x, ||\frac{\nabla(\mathcal{F}_l(t) - \mathcal{F}_{l_x}(t))}{\nabla \mathcal{F}^l(t)}|_{t=x}||_{-\infty} > c_l) > \alpha_l$, where $l \in [L]$ and $c_l, \alpha_l \in \mathbb{R}_+$.

$C_3$. $P_{x \sim D_x}(||v\frac{\nabla \mathcal{F}^l(t)}{\nabla t}|_{t=x}||_\infty \geq d_l||v||_\infty) > \beta_l$ for $\forall v \in \mathbb{R}^{1 \times n}$, where $l \in [L]$ and $d_l, \beta_l \in \mathbb{R}_+$.

$C_4$. $\{\mathcal{F}^l(x)\}_{x \in S}$ is in a low dimensional subspace of $\mathbb{R}^n$ and $||\mathcal{F}^l(x)||_0 > \gamma_l/n$, where $l \in [L]$, $x \in S$, and $\gamma_l \in \mathbb{R}_+$.

For $\gamma \in \mathbb{R}_+$, let $\mathcal{H}(\gamma)$ be the set of networks in $\mathcal{H}_\gamma(\Theta)$, whose accuracies are equal to or larger than that of $\mathcal{F}_\Theta$. Then

$$\min_{\widetilde{\mathcal{F}} \in H(\gamma)} \{R(\widetilde{\mathcal{F}}, D_x)\} \leq (1 - \rho)R(\mathcal{F}, D_x)$$

where $\rho = \frac{(\gamma c_1)^2 \alpha_1 \gamma_1 + \sum_{i=2}^L (\gamma c_i d_{i-1})^2 \gamma_i(\alpha_i + \beta_{i-1} - 1) + \beta_L(d_L \gamma)^2}{4A + (\gamma c_1)^2 \alpha_1 \gamma_1 + \sum_{i=2}^L (\gamma c_i d_{i-1})^2 \gamma_i(\alpha_i + \beta_{i-1} - 1) + \beta_L(d_L \gamma)^2}$. In other words, there exists an $\mathcal{F}_{\Theta_a}$ with adversarial rate $\geq \rho$.

We can make the robustness of the perturbed network as small as possible.

11

**Corollary 4.5.** *In Theorem 4.4, if $\alpha, \beta, c, d \in \mathbb{R}_+$ satisfy $\alpha_l > \alpha$, $\beta_l > \beta$, $c_l > c$, $d_l > d$, $\gamma_l > \gamma_{low}$ for $l \in [L]$, then*

$$\min_{\widetilde{\mathcal{F}} \in H(\gamma)} \{R(\widetilde{\mathcal{F}}, D_x)\} \le (1 - \frac{(\gamma c)^2 \alpha \gamma_{low} + (L-1)(\gamma c d)^2 \gamma_{low}(\alpha + \beta - 1) + \beta(d\gamma)^2}{4A + (\gamma c)^2 \alpha \gamma_{low} + (L-1)(\gamma c d)^2 \gamma_{low}(\alpha + \beta - 1) + \beta(d\gamma)^2}) R(\mathcal{F}, D_x).$$

*Furthermore, for $\rho \in (0,1)$, if $L > 1 + \frac{4(1-\rho)A}{\rho((\gamma c d)^2 \gamma_{low}(\alpha + \beta - 1))}$, then there exists an $\mathcal{F}_{\Theta_a} \in \mathcal{H}(\gamma)$ whose adversarial rate is $\ge 1 - \rho$.*

*Furthermore, for $\tau \in (0,1)$ satisfying $\tau < R(\mathcal{F}, D_x)$, if $L > \frac{4A(R(\mathcal{F}, D_x)/\tau - 1)}{(\gamma c d)^2 \gamma_{low}(\alpha + \beta - 1)} + 1$, then there exists an $\mathcal{F}_{\Theta_a} \in \mathcal{H}(\gamma)$ such that $R(\mathcal{F}_{\Theta_a}, D_x)\} \le \tau$.*

**Remark 4.4.** *From Corollary 4.5, if the depth of the DNN is sufficiently large, then there exist adversary parameters such that the attacked network has robustness measure as small as possible.*

**Remark 4.5.** *In practical computation, we use a finite set $T$ of samples satisfying $D_x$ and $R(\mathcal{F}, D_x)$ is approximately computed as $\widetilde{R}(\mathcal{F}, T) = 1/|T| \sum_{x \in T} R(\mathcal{F}, x)$. Since $\mathcal{F}_\Theta$ is a continued function in $\Theta$, if $\Theta_a$ is an adversarial parameter for $\Theta$ and $\widetilde{R}(\mathcal{F}, T)$ is used as the robustness measure, then there exists a small sphere $S_a$ with $\Theta_a$ as center such that all parameters in $S_a$ are also adversarial parameters for $\Theta$.*

The above remarks show that adversary parameters are inevitable in certain sense.

**Remark 4.6.** *In the model (18), two simplifications are made. However, the results proved in this section can be generalized to general DNNs. First, the bias vectors are not considered, which can be included as parts of the weight matrices by extending the input space slightly, similar to [22]. Second, it is assumed that $n_l = n$ for $l \in [L]$. This assumption could be removed by assuming $n = \max_{l \in [l]} n_i$.*

**Remark 4.7.** *Using $R_4(\mathcal{F}, D_x)$, results in theorem 4.4 cannot be obtained yet. But, we will use numerical experiments to show that the result is also valid for $R_4(\mathcal{F}, D_x)$.*

## 5 Experimental results

### 5.1 The setting

We use two networks: VGG19 [29] and Resnet56 [11]. We write VGG19 as $\mathcal{F}_V$ and Resnet56 as $\mathcal{F}_R$, which are trained with the adversarial training [20]. The experimental results are for the CIFAR-10 dataset.

We use both the adversary accuracy in (5) and the approximate robust radius in (6) to compute the adversarial rate. For a given data set $T$, the adversarial accuracy defined in (5) can be approximately computed with PGD [20] as follows

$$\widetilde{R}_3(\mathcal{F}_\Theta, T, \epsilon) = 1/|T| \sum_{x \in T} I(\widehat{\mathcal{F}}_\Theta(x') = l_x)$$

where $x' = \arg\max_{|\widetilde{x} - x|_\infty \le \epsilon} L_{\text{CE}}(\mathcal{F}_\Theta(\widetilde{x}), l_x)$. In the experiment, we set $\epsilon = 8/255$. The approximate robust radius in (6) can be computed as follows

$$\widetilde{R}_4(\mathcal{F}, T) = \frac{1}{T} \sum_{x \in T} \min_{l \ne l_x} \{ \frac{\mathcal{F}_{l_x}(x) - \mathcal{F}_l(x)}{||\nabla \mathcal{F}_{l_x}(x) - \nabla \mathcal{F}_l(x)||_1} I(\mathcal{F}_{l_x}(x_0) > \mathcal{F}_l(x_0)) \}$$

where the $L_1$-norm is used, since we consider $L_\infty$ adversarial samples.

The accuracies, adversarial accuracies, and AARs of $\mathcal{F}_V$ and $\mathcal{F}_R$ under the $L_\infty$ norm attack are given in Table 1, which are about the state of the art results for these DNNs.

| Net | AC | AA | AAR |
|---|---|---|---|
| $\mathcal{F}_V$ | 80% | 45% | 0.0770 |
| $\mathcal{F}_R$ | 83% | 49% | 0.0194 |

Table 1: Results for $\mathcal{F}_V$ and $\mathcal{F}_R$ on CIFAR-10. AC: accuracy, AA: adversarial accuracy, AAR: approximate accurate radius

## 5.2 Adversarial parameter attack

Let $\Theta$ be the parameter set of $\mathcal{F}_V$ or $\mathcal{F}_R$, and two kinds of parameter perturbation attacks will be carried out:

$L_{\infty,\gamma}$ **perturbation** for $\gamma \in \mathbb{R}_+$: We consider parameter perturbations in $B_\infty(\Theta, \Delta_\gamma)$, where $\Delta_\gamma = (\gamma|\theta_1|, \ldots, \gamma|\theta_k|)$ for $\Theta = (\theta_1, \ldots, \theta_k)$. In other words, $\gamma$ is the perturbation ratio. Also, the BN-layers will be changed to compute this kind of perturbations.

$L_{0,k}$ **perturbation** for $k \in \mathbb{N}_{>0}$: $k$ weight matrices are perturbed and $\max\{400, 1\%\#W_l\}$ pairs of weights are changed for $\mathcal{F}_V$ with the method given in section 3.2 ($1\%\#W_l$ pairs of weights are changed for $\mathcal{F}_R$), where $\#W_l$ is the number of entries of $W_l$. The BN-layers will not be changed to compute this kind of perturbations.

We set $\gamma_{\text{low}}$ in Remark 2.1 to be 90%, that is, if the accuracy of the perturbed network has is less than 90% of that of the original DNN, then the attack is considered failed.

### 5.2.1 Random parameter perturbation

We do random parameter perturbations and will use them as bases for comparisons. The results are given in Table 2.

| Attack | AC | AA | AR |
|---|---|---|---|
| No attack | 80% | 45% | 0 |
| $L_{\infty,0.02}$ | 80% | 39% | 0.13 |
| $L_{\infty,0.04}$ | 80% | 37% | 0.17 |
| $L_{\infty,0.06}$ | 79% | 36% | 0.2 |
| $L_{\infty,0.08}$ | 78% | 35% | 0.22 |
| $L_{\infty,0.10}$ | 78% | 34% | 0.24 |
| $L_{0,8}$ | 67% | 23% | 0.41(fail) |
| $L_{0,12}$ | 61% | 20% | 0.42(fail) |
| $L_{0,16}$ | 57% | 19% | 0.41(fail) |

Table 2: Random perturbations for $\mathcal{F}_V$ . AC: accuracy, AA: adversarial accuracy, AR: adversarial rate

For the $L_\infty$ perturbations, the accuracies are kept high, but the robustness does not decrease much, so the adversarial rates are low. For the $L_0$ perturbations, the accuracy decreases too much and are considered failed attacks. In either case, random perturbations are not good adversarial

13

parameters. Thus adversarial parameters are sparse around the trained parameters, which is similar with adversarial samples [40].

| Attack | AC | AA | AR |
|---|---|---|---|
| No attack | 83.1% | 48.5% | 0 |
| $L_{\infty,0.02}$ | 82.9% | 48.0% | 0.004 |
| $L_{\infty,0.04}$ | 82.7% | 48.0% | 0.011 |
| $L_{\infty,0.06}$ | 82.3% | 47.5% | 0.022 |
| $L_{\infty,0.08}$ | 81.7% | 46.6% | 0.039 |
| $L_{\infty,0.10}$ | 81.0% | 45.5% | 0.061 |
| $L_{0,10}$ | 81.5% | 47.8% | 0.016 |
| $L_{0,20}$ | 80.7% | 45.8% | 0.054 |
| $L_{0,30}$ | 80.5% | 45.7% | 0.057 |

Table 3: Random perturbations for $\mathcal{F}_R$. AC: accuracy, AA: adversarial accuracy, AR: adversarial rate

Results of random perturbations for $\mathcal{F}_R$ are given in Table 3. From the results, we can see that network $\mathcal{F}_R$ is much more robust against random parameter perturbations than $\mathcal{F}_V$.

### 5.2.2 Adversarial parameter attack on $\mathcal{F}_V$ and $\mathcal{F}_R$

We use algorithms in section 3 to create adversarial parameters. The training set $T$ contains 500 samples for which $\mathcal{F}$ give the correct label. The average results are given in Table 4.

| Attack | AC | AA in (5) | | ARR in (6) | |
|---|---|---|---|---|---|
| | | $\widetilde{R}_3(\mathcal{F},T)$ | AR | $\widetilde{R}_4(\mathcal{F},T)$ | AR |
| No attack | 80% | 45% | 0 | 0.0770 | 0 |
| $L_{\infty,0.02}$ | 78% | 38% | 0.15 | 0.0667 | 0.13 |
| $L_{\infty,0.04}$ | 77% | 30% | 0.32 | 0.0481 | 0.36 |
| $L_{\infty,0.06}$ | 76% | 22% | 0.49 | 0.0372 | 0.49 |
| $L_{\infty,0.08}$ | 76% | 10% | 0.74 | 0.0195 | 0.71 |
| $L_{\infty,0.10}$ | 77% | 8% | 0.79 | 0.0143 | 0.78 |
| $L_{0,8}$ | 72% | 27% | 0.36 | 0.0441 | 0.38 |
| $L_{0,12}$ | 76% | 24% | 0.44 | 0.0443 | 0.40 |
| $L_{0,16}$ | 74% | 22% | 0.47 | 0.0404 | 0.44 |

Table 4: Adversarial parameter attack for $\mathcal{F}_V$. AC: accuracy, AA: adversarial accuracy, AR: adversarial rate.

Comparing Tables 2 and 4, we can see that algorithms in section 3 can be used to create good adversarial parameters, especially for the $L_\infty$ attack. From Figure 1, we can see that the adversarial rate and adversarial accuracy for the $L_{\infty,\gamma}$ attack are near linear in $\gamma$ when $\gamma$ is small, and is gradually stabilized with the increase of $\gamma$. Also when $\gamma$ is very small, say $\gamma = 0.02$, the adversarial parameter attacks do not create good results, which means the network is approximately safe against these attacks for $\gamma \leq 0.02$. For $L_0$ attacks, we can see that the accuracies are increased lots comparing to the random perturbation and adversarial parameters are obtained successfully. Also, for the two kinds of robustness-measurements, the adversarial rates are very close. For network $\mathcal{F}_R$, similar results are obtained and are given in Table 5.
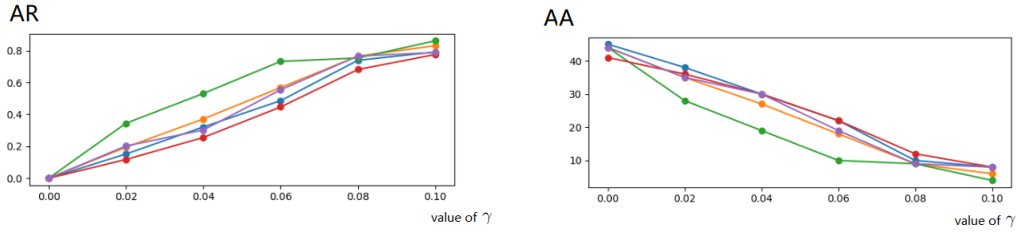
14

Figure 1: Left: Relation between $\gamma$ and the adversarial rate. Right: Relation between $\gamma$ and adversarial accuracy.

| Attack | AC | AA in (5) | | ARR in (6) | |
|---|---|---|---|---|---|
| | | $\widetilde{R}_3(\mathcal{F}, T)$ | AR | $\widetilde{R}_4(\mathcal{F}, T)$ | AR |
| No attack | 83% | 49% | 0 | 0.0194 | 0 |
| $L_{\infty, 0.02}$ | 84% | 39% | 0.20 | 0.0151 | 0.22 |
| $L_{\infty, 0.04}$ | 85% | 27% | 0.45 | 0.0127 | 0.35 |
| $L_{\infty, 0.06}$ | 86% | 14% | 0.72 | 0.0092 | 0.53 |
| $L_{\infty, 0.08}$ | 87% | 6% | 0.87 | 0.0064 | 0.67 |
| $L_{\infty, 0.10}$ | 87% | 1% | 0.98 | 0.0044 | 0.77 |
| $L_{0,10}$ | 80% | 26% | 0.45 | 0.0193 | 0 |
| $L_{0,20}$ | 78% | 18% | 0.59 | 0.0187 | 0.03 |
| $L_{0,30}$ | 72% | 9% | 0.71 | 0.0125 | 0.33 |

Table 5: Adversarial parameter attack for $\mathcal{F}_R$. AC: accuracy, AA: adversarial accuracy, AR: adversarial rate.

### 5.2.3 Affect of network depth and width on the adversarial parameter attack

We check how the network depth and width affect on the adversarial parameter attack. We use the $L_{\infty, \gamma}$ adversarial parameter attack for $\gamma = 0.02, 0.04$. Let $\mathcal{F}_V^k$ be the network which has the same width with $\mathcal{F}_V$ but has $k$ more layers, $\mathcal{F}_V(\alpha)$ the network which has the same depth with $\mathcal{F}_V$ but has $\alpha$ times width as $\mathcal{F}_V$. The results are given in Table 6. We can see that when the depth becomes larger, the attack becomes easier. This validates the results in section 4.2, for instance Corollary 4.5, where it shows that when the depth of the network becomes large, adversarial parameters exist.

The attack is much less sensitive to the width. The reason may be that there exist much redundancy on the width, similar to the results in [17, 25, 26], and the redundance can lead to limited search directions in the feature space and poor generalization performance, as shown in [21], so the attack is hard to improve when the width becomes larger.

| Network | $\gamma = 0$ | | $\gamma = 0.02$ | | | $\gamma = 0.04$ | | |
|---|---|---|---|---|---|---|---|---|
| | AC | AA | AC | AA | AR | AC | AA | AR |
| $\mathcal{F}_V$ | 80% | 45% | 78% | 38% | 0.15 | 77% | 30% | 0.32 |
| $\mathcal{F}_V^8$ | 80% | 47% | 76% | 37% | 0.20 | 78% | 32% | 0.31 |
| $\mathcal{F}_V^{16}$ | 78% | 44% | 74% | 35% | 0.19 | 75% | 27% | 0.37 |
| $\mathcal{F}_V^{24}$ | 79% | 43% | 73% | 30% | 0.28 | 73% | 20% | 0.49 |
| $\mathcal{F}_V^{32}$ | 76% | 44% | 72% | 28% | 0.34 | 71% | 19% | 0.53 |
| $\mathcal{F}_V(1.25)$ | 80% | 42% | 77% | 37% | 0.12 | 76% | 29% | 0.29 |
| $\mathcal{F}_V(1.5)$ | 81% | 41% | 78% | 36% | 0.12 | 77% | 30% | 0.26 |
| $\mathcal{F}_V(2)$ | 78% | 44% | 76% | 39% | 0.11 | 74% | 31% | 0.28 |
| $\mathcal{F}_V(2.5)$ | 80% | 44% | 79% | 35% | 0.20 | 76% | 30% | 0.30 |

Table 6: Affect of width and depth on adversarial parameter attack for $\mathcal{F}_V$. AC: accuracy, AA: adversarial accuracy, AR: adversarial rate.

We can use $\widetilde{R}_4(\mathcal{F}, T)$ to measure the robustness and similar results are obtained, which are given in Table 7.

| Network | $\gamma = 0$ | $\gamma = 0.02$ | | $\gamma = 0.04$ | |
|---|---|---|---|---|---|
| | $\widetilde{R}_4(\mathcal{F}, T)$ | $\widetilde{R}_4(\mathcal{F}, T)$ | AR | $\widetilde{R}_4(\mathcal{F}, T)$ | AR |
| $\mathcal{F}_V$ | 0.0770 | 0.0667 | 0.13 | 0.0481 | 0.36 |
| $\mathcal{F}_V^8$ | 0.0776 | 0.0686 | 0.11 | 0.0534 | 0.30 |
| $\mathcal{F}_V^{16}$ | 0.0815 | 0.0652 | 0.19 | 0.0479 | 0.40 |
| $\mathcal{F}_V^{24}$ | 0.0817 | 0.0630 | 0.21 | 0.0388 | 0.49 |
| $\mathcal{F}_V^{32}$ | 0.0808 | 0.0608 | 0.23 | 0.0392 | 0.48 |
| $\mathcal{F}_V(1.25)$ | 0.0750 | 0.0663 | 0.11 | 0.0524 | 0.29 |
| $\mathcal{F}_V(1.5)$ | 0.0763 | 0.0670 | 0.12 | 0.0563 | 0.25 |
| $\mathcal{F}_V(2)$ | 0.0750 | 0.0650 | 0.13 | 0.0499 | 0.32 |
| $\mathcal{F}_V(2.5)$ | 0.0775 | 0.0678 | 0.12 | 0.0489 | 0.35 |

Table 7: Affect of width and depth on adversarial parameter attack for $\mathcal{F}_V$. AR: adversarial rate

## 5.3 Direct adversarial parameters

We give experimental results for direct adversarial parameters introduced in section 2.2. We try to decrease the accuracies for samples with label 0 and keep the accuracies and robustness for other samples. The experimental results are for the network $\mathcal{F}_V$ and CIFAR-10 and are given in Table 8.

| Attack | $AC_1$ | $AA_1$ | $AC_0$ | AR |
|---|---|---|---|---|
| $L_{\infty,0.02}$ | 77% | 35% | 11% | 0.65 |
| $L_{\infty,0.04}$ | 78% | 40% | 3% | 0.83 |
| $L_{\infty,0.06}$ | 79% | 42% | 1% | 0.92 |
| $L_{\infty,0.08}$ | 80% | 43% | 1% | 0.95 |
| $L_{\infty,0.1}$ | 80% | 45% | 1% | 0.99 |

Table 8: Direct adversarial parameter attack for $\mathcal{F}_V$. $AC_1$ and $AA_1$ are for samples with label $\neq 0$, $AC_0$ is the accuracy for samples with label 0.

Comparing to Tables 8 and 4, we can see that direct adversarial parameters for a given label are much easier to compute than adversarial parameters. High quality direct adversarial parameters

can be obtained by using perturbation ratios $6\% - 10\%$. Results for network $\mathcal{F}_R$ are given in Table 9, from which we can see that it is slightly more difficult to attack $\mathcal{F}_R$.

| Attack | $AC_1$ | $AA_1$ | $AC_0$ | AR |
|---|---|---|---|---|
| $L_{\infty,0.02}$ | 82% | 46% | 52% | 0.38 |
| $L_{\infty,0.04}$ | 81% | 47% | 35% | 0.57 |
| $L_{\infty,0.06}$ | 81% | 47% | 10% | 0.83 |
| $L_{\infty,0.08}$ | 81% | 46% | 4% | 0.90 |
| $L_{\infty,0.1}$ | 80% | 46% | 1% | 0.93 |

Table 9: Direc adversarial parameter attack for $\mathcal{F}_R$. $AC_1$ and $AA_1$ are for samples with label $\neq 0$, $AC_0$ is the accuracy for samples with label 0.

## 5.4 Adversarial parameters for a given sample

We give experimental results for adversarial parameters for a given sample introduced in section 2.2. $\widetilde{R}_2(\mathcal{F},x) = \min_{i \neq l_x}\{\frac{\mathcal{F}_{l_x}(x)-\mathcal{F}_i(x)}{||\nabla\mathcal{F}_{l_x}(x)-\nabla\mathcal{F}_i(x)||_1}\}$ is used to measure the robustness of $\mathcal{F}$ at sample $x$. Let $S$ be a subset of the test set containing 100 samples such that $\mathcal{F}$ gives the correct label for all of them and all samples in $S$ are robust in that, no adversarial samples were found using PGD-10 with $L_\infty = \frac{8}{255}$.

| Attack | $\widetilde{R}_2(\mathcal{F},x)$ | $N_1$ | $N_2$ | AR |
|---|---|---|---|---|
| before attack | 0.078 | 100 | 100 | 0 |
| $L_{\infty,0.02}$ | 0.016 | 0 | 100 | 0.79 |
| $L_{\infty,0.04}$ | 0.010 | 0 | 100 | 0.87 |
| $L_{\infty,0.06}$ | 0.008 | 0 | 100 | 0.89 |
| $L_{\infty,0.08}$ | 0.006 | 0 | 100 | 0.92 |
| $L_{\infty,0.1}$ | 0.005 | 0 | 100 | 0.94 |

Table 10: Adversarial parameter attack to $\mathcal{F}_V$ for a given sample. AR: adversarial rate

For each sample $x \in S$, we compute $L_{\infty,\gamma}$ adversarial parameters and the average results are given in Table 10, where $N_1$ is the number of robust samples, and $N_2$ the number of samples which are given the correct labels. Comparing to Tables 8, 4, and 10, we can see that adversarial rates for a single sample are about the same as that for a given label. Similar results for network $\mathcal{F}_R$ are given in Table 11.

| Attack | $\widetilde{R}_2(\mathcal{F},x)$ | $N_1$ | $N_2$ | AR |
|---|---|---|---|---|
| before attack | 0.057 | 100 | 100 | 0 |
| $L_{\infty,0.02}$ | 0.008 | 0 | 100 | 0.86 |
| $L_{\infty,0.04}$ | 0.008 | 0 | 100 | 0.86 |
| $L_{\infty,0.06}$ | 0.008 | 0 | 100 | 0.86 |
| $L_{\infty,0.08}$ | 0.008 | 0 | 100 | 0.86 |
| $L_{\infty,0.1}$ | 0.008 | 0 | 100 | 0.86 |

Table 11: Adversarial parameter attack to $\mathcal{F}_R$ for a given sample. AR: Adversarial Rate

17

# 6 Proofs for the theorems in section 4

## 6.1 Proofs of Theorems 4.1 and 4.2

We introduce several notations. Let $||x||_{-\infty} = \min_{i\in[n]}\{|x_i|\}$ for $x \in \mathbb{R}^n$, and $||W||_{-\infty,2} = \min_{i\in[a]}\{||W^{(i)}||_2\}$ for $W \in \mathbb{R}^{a\times b}$, where $W^{(i)}$ is the $i$-th row of $W$. If $\mathcal{F}$ is a network, we use $\mathcal{F}_i(x)$ to denote the $i$-th coordinate of $\mathcal{F}(x)$. A lemma is proved first.

**Lemma 6.1.** *Let $v \in \mathbb{R}^n$ and $v \neq 0$. Then there exists a vector $w \in \mathbb{R}^n$ such that $w \perp v$, $||w||_\infty = 1$ and $||w||_2 \geq \sqrt{n-1}$.*

*Proof.* Let $S = \arg\min_{S\subseteq[n]}\{|\sum_{i\in S}|v_i| - \sum_{i\in[n]\setminus S}|v_i||\}$. We can assume $\sum_{i\in S}|v_i| - \sum_{i\in[n]\setminus S}|v_i| = k \geq 0$. For any $j \in S$ such that $v_j \neq 0$ and $S_1 = S/\{j\}$, we have $|\sum_{i\in S_1}|v_i| - \sum_{i\in[n]/S_1}|v_i|| = |2|v_j| - k| \geq k$, which means $k \leq |v_j|$.

We now define $w \in \mathbb{R}^n$. Set $w_i = 1$ if $v_i = 0$. Select a $j \in S$ such that $v_j \neq 0$ and let $w_i = \text{sign}(v_i)$ if $i \in S/\{j\}$ and $v_i \neq 0$, and $w_j = \frac{-\sum_{i\in S}|v_i|+\sum_{i\in[n]\setminus S}|v_i|+|v_j|}{v_j}$. For $i \in [n] \setminus S$, let $w_i = -\text{sign}(v_i)$ if $v_i \neq 0$. It is easy to check that $||w||_\infty = 1$, $||w||_2 \geq \sqrt{n-1}$, and $w \perp v$. The lemma is proved. $\square$

We now prove Theorem 4.1.

*Proof.* By Lemma 6.1, there exists a vector $v \in \mathbb{R}^n$ such that $v \perp x_0$, $||v||_2 \geq \sqrt{n-1}$ and $||v||_\infty = 1$. Moreover, we can assume that at least $\eta n_1/2$ coordinates of $\text{Relu}(W_1(x+\epsilon v)+b_1)$ are bigger than $b$. If this is not valid, we just need to change $v$ to $-v$, and then $\text{Relu}(W_1(x + \epsilon v) + b_1) + \text{Relu}(W_1(x - \epsilon v)+b_1) \geq 2\text{Relu}(W_1 x + b_1)$, since $\text{Relu}(x) + \text{Relu}(y) \geq \text{Relu}(x+y)$ for all $x, y \in \mathbb{R}$. By condition $C_3$, at least $\eta n_1$ coordinates of $2\text{Relu}(W_1 x + b_1)$ are bigger than $2b$, but fewer than $\eta n_1/2$ coordinates of $\text{Relu}(W_1(x + \epsilon v) + b_1)$ are bigger than $b$, so at least $\eta n_1/2$ coordinates of $\text{Relu}(W_1(x - \epsilon v) + b_1)$ are bigger than $b$.

Let $l_2 \in [m]$ such that $l_2 \neq l_{x_0}$, $\overline{W}_2 = -(W_2^{(l_{x_0})} - W_2^{(l_2)}) \in \mathbb{R}^{1\times n_1}$, $U_v \in \mathbb{R}^{n_1\times n}$ all of whose rows are $\gamma v^\tau$ (the transposition of $v$), and $U = \text{diag}(\text{sign}(\overline{W}_2)) \in \mathbb{R}^{n_1\times n_1}$. Let $\widetilde{W}_1 = W_1 + UU_v$ and

$$\widetilde{\mathcal{F}}(x) = W_2\text{Relu}(\widetilde{W}_1 x + b_1) + b_2.$$

We will show that $\widetilde{\mathcal{F}}$ satisfies the condition of the theorem.

Since $v \perp x_0$, we have $\widetilde{\mathcal{F}}(x_0) = \mathcal{F}(x_0)$ and $\widetilde{\mathcal{F}}$ gives the correct label for $x_0$. Since $||v||_\infty = 1$, we have $||\widetilde{W}_1 - W_1||_\infty = ||UU_v||_\infty = ||U_v||_\infty = ||\gamma v||_\infty \leq \gamma$, and thus $\widetilde{\mathcal{F}}(x) \in H_\gamma(\Theta)$.

So it suffices to show that $\widetilde{\mathcal{F}}(x_0 + \epsilon v)$ will not give $x_0 + \epsilon v$ label $l_{x_0}$, which means that $\widetilde{\mathcal{F}}$ has adversarial samples to $x_0$ in $S_\infty(x_0, \epsilon)$. Since

$$\begin{aligned}
&\widetilde{\mathcal{F}}(x_0 + \epsilon v)\\
=\ & W_2\text{Relu}(\widetilde{W}_1(x_0 + \epsilon v) + b_1) + b_2\\
=\ & W_2\text{Relu}(W_1(x_0 + \epsilon v) + b_1 + \epsilon UU_v v) + b_2
\end{aligned}$$

we have

$$\begin{aligned}
&\widetilde{\mathcal{F}}_{l_{x_0}}(x_0 + \epsilon v) - \widetilde{\mathcal{F}}_{l_2}(x_0 + \epsilon v)\\
=\ & \mathcal{F}_{l_{x_0}}(x_0 + \epsilon v) - \mathcal{F}_{l_2}(x_0 + \epsilon v)+\\
& \overline{W}_2(\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v)).
\end{aligned}$$

18

Since $e(\text{Relu}(f) - \text{Relu}(e + f)) = -|f|(|\text{Relu}(e) - \text{Relu}(e + f)|)$ for all $e, f \in \mathbb{R}$, we have

$$-\overline{W}_2(\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v))$$
$$= |\overline{W}_2|(|\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v)|).$$

Since $\epsilon UU_v v = \epsilon \gamma ||v||_2^2 \text{sign}(\overline{W}_2)$ and $||v||_2 \geq n-1$, each weight of $|\epsilon UU_v v|$ is at least $\epsilon \gamma (n-1)$. Note that if $e > 0$ and $f \in \mathbb{R}$, then $|\text{Relu}(e) - \text{Relu}(e - f)| \geq \min\{e, |f|\}$. As a consequence, if $i$ satisfies $(\text{Relu}(W_1(x + \epsilon v) + b_1))_i > b$, then $(|\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v)|)_i \geq \min\{\epsilon \gamma (n-1), b\}$. Since at least $\eta n_1/2$ coordinates of $\text{Relu}(W_1(x + \epsilon v) + b_1)$ are bigger than $b$, we have $||\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v)||_1 \geq \eta n_1/2 \min\{\epsilon \gamma (n-1), b\}$. By condition $C_2$, it is easy see

$$-\overline{W}_2(\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v))$$
$$= |\overline{W}_2|(|\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v)|)$$
$$\geq ||\overline{W}_2||_{-\infty}||\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v)||_1$$
$$\geq \min\{\epsilon \gamma (n-1), b\} c n_1 \eta/2,$$

that is, $\overline{W}_2(\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v)) \leq -\min\{\epsilon \gamma (n-1), b\} c n_1 \eta/2$. By condition $C_1$, we have $\mathcal{F}_{l_{x_0}}(x_0 + \epsilon v) - \mathcal{F}_{l_2}(x_0 + \epsilon v) \leq A$. Then we have

$$\widetilde{\mathcal{F}}_{l_{x_0}}(x_0 + \epsilon v) - \widetilde{\mathcal{F}}_{l_2}(x_0 + \epsilon v)$$
$$= \mathcal{F}_{l_{x_0}}(x_0 + \epsilon v) - \mathcal{F}_{l_2}(x_0 + \epsilon v) +$$
$$\overline{W}_2(\text{Relu}(W_1(x + \epsilon v) + b_1) - \text{Relu}(W_1(x + \epsilon v) + b_1 + \epsilon UU_v v))$$
$$\leq A - \min\{\epsilon \gamma (n-1), b\} c n_1 \eta/2 < 0.$$

Thus if $n_1 > \frac{2A}{\min\{\epsilon \gamma (n-1), b\} c \eta}$, then $\widetilde{\mathcal{F}}_{l_{x_0}}(x_0 + \epsilon v) - \widetilde{\mathcal{F}}_{l_2}(x_0 + \epsilon v) < 0$ and the label of $\widetilde{\mathcal{F}}(x_0 + \epsilon v)$ is not $l_{x_0}$. The theorem is proved. $\square$

We now prove Theorem 4.2.

*Proof.* By condition $C_4$, for $l \in [m]$, there exist $v_l \in \mathbb{R}^n$ such that $v_l \perp S$, $v_l \perp v_k$ for $l \neq k$, $||v_l||_2 = 1$. Then $||v_l||_\infty \leq 1$.

By condition $C_3$, at least $\eta n_1/2$ coordinates of $\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1)$ are bigger than $b$ or at least $\eta n_1/2$ coordinates of $\text{Relu}(W_1(x - \epsilon v_{l_x}) + b_1)$ are bigger than $b$, similar to the proof of Theorem 4.1.

For convenience, we write $G(x, y) : (\mathbb{R}^n, \mathbb{R}) \to \mathbb{R}$ as $G(x, y) = ||\text{sign}(x - yI_n)||_0$, where $I_n$ is the vector with entries 1. It is easy to see that, $G(x, b)$ is the number of coordinates of $x$ that are bigger than $b$. So we have $G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2$ or $G(\text{Relu}(W_1(x - \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2$ for all $x$, and hence for $l \in [m]$, we have

$$P_{x \sim D_x}(G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2 \text{ or } G(\text{Relu}(W_1(x - \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2 \,|\, l_x = l) = 1.$$

For events $e$ and $f$, $P(e \text{ or } f) \leq P(e) + P(f)$. We thus have $P_{x \sim D_x}(G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2 \,|\, l_x = l) \geq 0.5$ or $P_{x \sim D_x}(G(\text{Relu}(W_1(x - \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2 \,|\, l_x = l) \geq 0.5$. Without loss of generality, we can assume that for any $l \in [m]$, $P_{x \sim D_x}(G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2 \,|\, l_x = l) \geq 0.5$. Therefore,

$$P_{x \sim D_x}(G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2)$$
$$= \sum_{l \in [m]} P_{x \sim D_x}(l_x = l) P_{x \sim D_x}(G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) \geq \eta n_1/2 \,|\, l_x = l)$$
$$\geq 0.5 \sum_{l \in [m]} P_{x \sim D_x}(l_x = l)$$
$$= 0.5.$$

For $l \in [m]$, let $\overline{W}_2^{(l)} = W_2^{(l)} - W_2^{(l+1)}$, where $W_2^{(m+1)} = W_2^{(1)}$. Now assume $U_v^l \in \mathbb{R}^{n_1 \times n}$, whose rows are all $\gamma v_l$, and $U_l = \text{diag}(\text{sign}(\overline{W}_2^{(l)}))$. Let $\widetilde{W}_1 = W_1 + \frac{1}{m} \sum_{l=1}^{m} U_l U_v^l$, and

$$\widetilde{\mathcal{F}}(x) = W_2 \text{Relu}(\widetilde{W}_1 x + b_1) + b_2.$$

We will show that $\widetilde{\mathcal{F}}(x)$ satisfies the conditions of the theorem.

It is easy to see that $\widetilde{\mathcal{F}}$ is in $\mathcal{H}_\gamma(\Theta)$. For any $x \in S$, $\widetilde{W}_1 x = W_1 x + \frac{1}{m} \sum_{l=1}^{m} (U_l U_v^l) x = W_1 x$, which means $\widetilde{\mathcal{F}}(x) = \mathcal{F}(x)$ and the accuracy of $\widetilde{\mathcal{F}}$ over $D_x$ is equal to that of $\mathcal{F}$.

Let $x \in S$ satisfy $G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) > \eta n_1/2$ and $l_2 \neq l_x$. By conditions $C_1$ and $C_2$ and similar to the proof of Theorem 4.1, we have

$$
\begin{aligned}
& \widetilde{\mathcal{F}}_{l_x}(x + \epsilon v_{l_x}) - \widetilde{\mathcal{F}}_{l_2}(x + \epsilon v_{l_x}) \\
= \ & \mathcal{F}_{l_x}(x + \epsilon v_{l_x}) - \mathcal{F}_{l_2}(x + \epsilon v_{l_x}) + \\
& W_c^{(l_x)}(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1) - \text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1 + \epsilon U_l U_v^l v/m)) \\
\leq \ & A - \min\{\epsilon \gamma/m, b\} c n_1 \eta/2 < 0.
\end{aligned}
$$

Thus $\widetilde{\mathcal{F}}$ does not give label $l_x$ to $x + \epsilon v_{l_x}$ and $\widetilde{\mathcal{F}}$ has an adversarial sample of $x$ in $S_\infty(x, \epsilon)$. Furthermore, since $P_{x \sim D_x}(G(\text{Relu}(W_1(x + \epsilon v_{l_x}) + b_1), b) \geq \gamma n_1/2) > 0.5$, we have

$$P_{x \sim D_x}(\widetilde{\mathcal{F}} \text{ has an adversarial sample of } x \text{ in } S_\infty(x, \epsilon)) > 0.5.$$

The theorem is proved. $\qquad \square$

## 6.2 Proofs of Theorems 4.3 and 4.4

We first prove two lemmas.

**Lemma 6.2.** *For $l \in [m]$, let $S_l$ be a non-empty bounded closed subset of $\mathbb{R}^{n \times n}$ such that $W \in S_l$ implies $-W \in S_l$. Also let $S_0$ be a non-empty bounded closed subset of $\mathbb{R}^{1 \times n}$ such that $x \in S_0$ implies $-x \in S_0$. Let $U_0 \in \mathbb{R}^{1 \times n}$ and $U_l \in \mathbb{R}^{n \times n}$ for $l \in [m]$. Define maps: $T_0(x) : \mathbb{R}^{1 \times n} \to \mathbb{R}^{1 \times n}$ by $T_0(x) = x \prod_{l=1}^{m} U_l$, and for $l \in [m]$, $T_l(W) : \mathbb{R}^{n \times n} \to \mathbb{R}^{1 \times n}$ by $T_l(W) = U_0(\prod_{j=1}^{l-1} U_j) W (\prod_{j=l+1}^{m} U_j)$. Then*

$$\max_{x_l \in S_l, \forall 0 \leq l \leq m} \{\|\prod_{l=0}^{m}(x_l + U_l)\|_2^2\} \geq \|\prod_{l=0}^{m} U_l\|_2^2 + \sum_{l=0}^{m} \max_{x_l \in S_l}\{\|T_l(x_l)\|_2^2\}$$

*Proof.* For $l \in [m]$, let $u_l = \arg\max_{x \in S_l} \|T_l(x)\|_2^2$, which exists because $S_l$ is bounded and closed. Then

$$
\begin{aligned}
& \max_{x_l \in S_l, \forall 0 \leq l \leq m}\{\|\prod_{l=0}^{m}(x_l + U_l)\|_2^2\} \\
\geq \ & \max_{x_l \in \{u_l, -u_l\}, \forall 0 \leq l \leq m}\{\|\prod_{l=0}^{m}(x_l + U_l)\|_2^2\} \\
= \ & \frac{1}{2^{m+1}} \sum_{x_l \in \{u_l, -u_l\}, \forall 0 \leq l \leq m} \|\prod_{l=0}^{m}(x_l + U_l)\|_2^2 \\
= \ & \sum_{M_l \in \{u_l, U_l\}, \forall 0 \leq l \leq m} \|\prod_{l=0}^{m} M_l\|_2^2 \\
\geq \ & \|\prod_{l=0}^{m} U_l\|_2 + \sum_{l=0}^{m} \|T_l(u_l)\|_2^2 \\
= \ & \|\prod_{l=0}^{m} U_l\|_2 + \sum_{l=0}^{m} \max_{x_l \in S_l}\{\|T_l(x_l)\|_2^2\}
\end{aligned}
$$

The lemma is proved. $\qquad \square$

**Lemma 6.3.** *For $l \in [m]$, let $S_l$ be a non-empty closed subset of bounded functions from $\mathbb{I}^k$ to $R^{n \times n}$ such that $c(x) \in S_l$ implies $-c(x) \in S_l$. Also let $S_0$ be a non-empty closed subset of bounded functions from $\mathbb{I}^k$ to $\mathbb{R}^{1 \times n}$ such that $c(x) \in S_0$ implies $-c(x) \in S_0$. Assume $U_0(x) : \mathbb{I}^k \to \mathbb{R}^{1 \times n}$, $U_l(x) : \mathbb{I}^k \to$*

$\mathbb{R}^{n \times n}$ for $l \in [m]$. Define maps: $T_0(c, x) : (S_0, \mathbb{I}^k) \to \mathbb{R}^{1 \times n}$ by $T_0(c, x) = c(x) \prod_{l=1}^{m} U_l(x)$; and $T_l(c, x) : (S_l, \mathbb{I}^k) \to \mathbb{R}^{1 \times n}$ by $T_l(x) = (\prod_{j=0}^{l-1} U_j(x)) c(x) (\prod_{j=l+1}^{m} U_j(x))$ for $l \in [m]$. Then we have

$$\max_{c_l(x) \in S_l, \forall 0 \le l \le m} \{ \int_{x \sim D_x} || \prod_{l=0}^{m} (c_l(x) + U_l(x)) ||_2^2 \} \ge \int_{x \sim D_x} || \prod_{l=0}^{m} U_l(x) ||_2^2 + \sum_{l=0}^{m} \max_{c_l(x) \in S_l} \{ \int_{x \sim D_x} ||T_l(c_l, x)||_2^2 \}.$$

*Proof.* Denote $\overline{c_l}(x) = \arg\max_{c \in S_l} || \int_{x \sim D_x} T_l(c, x) ||_2^2$, which must exist because $S_l$ is closed and its elements are bounded. Then

$$
\begin{aligned}
& \max_{c_l(x) \in S_l, \forall 0 \le l \le m} \{ \int_{x \sim D_x} || \prod_{l=0}^{m} (c_l(x) + U_l(x)) ||_2^2 \} \\
\ge\ & \max_{c_l(x) \in \{\overline{c_l}(x), -\overline{c_l}(x)\}, \forall 0 \le l \le m} \{ \int_{x \sim D_x} || \prod_{l=0}^{m} (c_l(x) + U_l(x)) ||_2^2 \} \\
\ge\ & \frac{1}{2^{m+1}} \sum_{c_l \in \{\overline{c_l}(x), -\overline{c_l}(x)\}, \forall 0 \le l \le m} \int_{x \sim D_x} || \prod_{l=0}^{m} (c_l(x) + U_l(x)) ||_2^2 \\
=\ & \sum_{M_l(x) \in \{\overline{c_l}(x), U_l(x)\}, \forall 0 \le l \le m} || \int_{x \sim D_x} \prod_{l=0}^{m} M_l(x) ||_2^2 \\
\ge\ & \int_{x \sim D_x} || \prod_{l=0}^{m} U_l(x) ||_2^2 + \sum_{l=0}^{m} \int_{x \sim D_x} ||T_l(\overline{c_l}, x)||_2^2 \\
=\ & \int_{x \sim D_x} || \prod_{l=0}^{m} U_l ||_2 + \sum_{l=0}^{m} \max_{c_l \in S_l} \{ \int_{x \sim D_x} ||T_l(c_l, x)||_2^2 \}.
\end{aligned}
$$

The lemma is proved. $\qquad\square$

We now prove Theorem 4.3.

*Proof.* Let $S_l$ be the set of $Q \in \mathbb{R}^{n \times n}$ satisfying $||Q||_\infty \le \gamma$ and $Q\mathcal{F}^{l-1}(x_0) = 0$ for $l \in [L+1]$. Note that $Q$ satisfies $n(l+1)$ linear equations and has $n^2$ variables. Since $n \gg l$ in DNNs, we can assume that $S_l$ is not empty. Let $\mathbb{F}$ be the set of networks $\widetilde{\mathcal{F}} : \mathbb{R}^n \to \mathbb{R}^m$ satisfying

$$\widetilde{\mathcal{F}}(x) = \widetilde{W}_{L+1} \sigma(\widetilde{W}_L \sigma(\widetilde{W}_{L-1} \sigma(\ldots \sigma(\widetilde{W}_1 x))))$$

where $\widetilde{W}_l - W_l \in S_l$ for all $l \in [L+1]$. We have $\widetilde{\mathcal{F}}^l(x_0) = \mathcal{F}^l(x_0)$ for $l \in [L+1]$. So $\widetilde{\mathcal{F}}(x_0) = l_{x_0}$ if $\widetilde{\mathcal{F}} \in \mathbb{F}$. It is also easy to see that $\mathbb{F} \subset \mathcal{H}_\gamma(\Theta)$. So it suffices to prove

$$\min_{\widetilde{\mathcal{F}} \in \mathbb{F}} \{ \frac{R(\widetilde{\mathcal{F}}, x_0)}{R(\mathcal{F}, x_0)} \} \le 1 - \frac{\gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)}{4A + \gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)}.$$

Let $l_2 = \arg\min_{i \ne l_{x_0}} \{ \frac{|\mathcal{F}_{l_{x_0}}(x_0) - \mathcal{F}_i(x_0)|^2}{||\nabla(\mathcal{F}_{l_{x_0}}(x_0)) - \nabla(\mathcal{F}_i(x_0))||_2^2} I(\mathcal{F}_{l_{x_0}}(x_0) > \mathcal{F}_i(x_0)) \}$. Then

$$R(\mathcal{F}, x_0) = \frac{|\mathcal{F}_{l_{x_0}}(x_0) - \mathcal{F}_{l_2}(x_0)|^2}{||\nabla(\mathcal{F}_{l_{x_0}}(x_0)) - \nabla(\mathcal{F}_{l_2}(x_0))||_2^2}.$$

Then for all $\widetilde{\mathcal{F}} \in \mathbb{F}$, we have

$$R(\widetilde{\mathcal{F}}, x_0) \le \frac{|\widetilde{\mathcal{F}}_{l_x}(x_0) - \widetilde{\mathcal{F}}_{l_2}(x_0)|^2}{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||_2^2} = \frac{|\mathcal{F}_{l_x}(x_0) - \mathcal{F}_{l_2}(x_0)|^2}{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||_2^2}$$

So we have

$$
\begin{aligned}
& \min_{\widetilde{\mathcal{F}} \in \mathbb{F}} \{ \frac{R(\widetilde{\mathcal{F}}, x_0)}{R(\mathcal{F}, x_0)} \} \\
\le\ & \min_{\widetilde{\mathcal{F}} \in \mathbb{F}} \{ \frac{|\mathcal{F}_{l_x}(x_0) - \mathcal{F}_{l_2}(x_0)|^2}{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||_2^2} / \frac{|\mathcal{F}_{l_x}(x) - \mathcal{F}_{l_2}(x)|^2}{||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2} \} \qquad (19) \\
\le\ & \min_{\widetilde{\mathcal{F}} \in \mathbb{F}} \{ \frac{||\nabla(\mathcal{F}_{l_x}(x_0)) - \nabla(\mathcal{F}_{l_2}(x_0))||_2^2}{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||_2^2} \}.
\end{aligned}
$$

To prove the theorem, we will first find a lower bound for $\max_{\widetilde{\mathcal{F}}\in\mathbb{F}}\{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0))-\nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||_2\}$. Let $J_l(x)=\text{diag}(\text{sign}(\mathcal{F}^l(x))): \mathbb{R}^n\to\mathbb{R}^{n\times n}$ for $l\in[L]$. Then

$$\frac{\nabla\mathcal{F}_i(x)}{\nabla x}=W^i_{L+1}(J_L(x)W_L)(J_{L-1}(x)W_{L-1})\ldots(J_1(x)W_1).$$

Let $A_l(x)=\frac{\nabla\mathcal{F}^L(t)}{\nabla\mathcal{F}^l(t)}|_{t=x}$ and $B_l(x)=\frac{\nabla\mathcal{F}^l(t)}{\nabla t}|_{t=x}$ for $l\in[L]$. Then we have

$$A_l(x)=(J_L(x)W_L)(J_{L-1}(x)W_{L-1})\ldots(J_{l+1}(x)W_{l+1})$$
$$B_l(x)=(J_l(x)W_l)(J_{l-1}(x)W_{l-1})\ldots(J_1(x)W_1).$$

Let $A_l=A_l(x_0)$, $B_l=B_l(x_0)$, $J_l=J_l(x_0)$. Then for all $\widetilde{\mathcal{F}}\in\mathbb{F}$,

$$\frac{\nabla\widetilde{\mathcal{F}}_i(x)}{\nabla x}|_{x=x_0}=\widetilde{W}^i_{L+1}(J_L\widetilde{W}_L)(J_{L-1}\widetilde{W}_{L-1})\ldots(J_1\widetilde{W}_1).$$

Denote $\overline{W}_l=\widetilde{W}_l-W_l\in S_l$ for $l\in[L+1]$. We have

$$\begin{aligned}
&\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0))-\nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))\\
=&(\widetilde{W}^{(l_x)}_{L+1}-\widetilde{W}^{(l_2)}_{L+1})(J_L\widetilde{W}_L)(J_{L-1}\widetilde{W}_{L-1})\ldots(J_1\widetilde{W}_1)\\
=&(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1}+\overline{W}^{(l_x)}_{L+1}-\overline{W}^{(l_2)}_{L+1})(J_L(W_L+\overline{W}_L))(J_{L-1}(W_{L-1}+\overline{W}_{L-1}))\ldots(J_1(W_1+\overline{W}_1)).
\end{aligned}$$

Now, let $\gamma(\overline{W}_1)=(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_1J_1\overline{W}_1$, $M_i(\overline{W}_i)=(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_iJ_i\overline{W}_iB_{i-1}$, where $i\in\{2,3,\ldots,L\}$, $M_{L+1}(\overline{W}_{L+1})=(\overline{W}^{(l_x)}_{L+1}-\overline{W}^{(l_2)}_{L+1})B_L$. By Lemma 6.2, we have

$$\begin{aligned}
&\max_{\widetilde{\mathcal{F}}\in\mathbb{F}}\{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0))-\nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||^2_2\}\\
\geq&\max_{\overline{W}_j\in S_j,j\in[L+1]}\{||\nabla(\mathcal{F}_{l_x}(x_0))-\nabla(\mathcal{F}_{l_2}(x_0))||^2_2+\textstyle\sum_{i=1}^{L+1}||M_i(\overline{W}_i)||^2_2\}\\
=&||\nabla(\mathcal{F}_{l_x}(x_0))-\nabla(\mathcal{F}_{l_2}(x_0))||^2_2+\textstyle\sum_{i=1}^{L+1}\max_{\overline{W}_i\in S_i}\{||M_i(\overline{W}_i)||^2_2\}.
\end{aligned}$$

It is easy to see for any $l\in[L]$, $(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_l=\frac{\nabla\mathcal{F}_{l_x}(t)-\mathcal{F}_{l_2}(t)}{\nabla\mathcal{F}^l(t)}|_{t=x_0}$, so by condition $C_3$, we have $||(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_i||_{-\infty}>c$. Let $S^1_l$ be the subset of $S_l$ containing those $C$ which have at most one nonzero row. Hence, for $x\in\mathbb{R}^{1\times n}$ and $M\in\mathbb{R}^{n\times n}$, if at most one row of $M$ is nonzero, we have $||xM||_\infty=\max_{i,j\in[n]}\{|x_iM_{i,j}|\}\geq||x||_{-\infty}||M||_\infty$, where $x_i$ is the $i$-th weight of $x$, $M_{i,j}$ is the weight of $M$ at $i$-th row and $j$-th column. Thus

$$\begin{aligned}
&\max_{\overline{W}_1\in S_1}\{||\gamma(\overline{W}_1)||_2\}\\
=&\max_{\overline{W}_1\in S_1}\{||(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_1J_1\overline{W}_1||_2\}\\
\geq&\max_{\overline{W}_1\in S_1}\{||(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_1J_1\overline{W}_1||_\infty\}\\
\geq&\max_{\overline{W}_1\in S^1_1}\{||(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_1J_1\overline{W}_1||_\infty\}\\
\geq&||(W^{(l_x)}_{L+1}-W^{(l_2)}_{L+1})A_1||_{-\infty}\max_{\overline{W}_1\in S^1_1}\{||J_1\overline{W}_1||_\infty\}\\
\geq&\gamma c.
\end{aligned}$$

Moreover, by condition $C_4$, there exists a column $L_{i-1}$ of $B_{i-1}$ such that $\pi-r\geq\alpha(\mathcal{F}^{i-1}(x_0),L_{i-1})\geq r$, where $\alpha(x,y)$ is the angle between $x,y$. Therefore, there exists a vector $v_i\in\mathbb{R}^n$ such that

22

$v \perp \mathcal{F}^{i-1}(x)$, $||v_i||_\infty = \gamma$ and consider condition $C_2$, we have $\langle v_i, L_{i-1} \rangle = ||v_i||_2 ||L_{i-1}||_2 \cos(\pi/2 - r) \geq \sin(r)b\gamma$.

Then $\max_{\overline{W}_i \in S_i^1} ||J_i \overline{W}_i B_{i-1}||_\infty \geq \sin(r)b\gamma$, because there must exist a $\overline{W}_i \in S_i^1$ whose only nonzero row is $v_i$ and $J_i \overline{W}_i = \overline{W}_i$. For $l \in \{2, 3, \ldots, L\}$, by condition $C_3$, we have

$$
\begin{aligned}
& \max_{\overline{W}_l \in S_l} \{||M_l(\overline{W}_l)||_2 \\
=\ & \max_{\overline{W}_l \in S_l} \{||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_i J_l \overline{W}_l B_{l-1}||_2\} \\
\geq\ & \max_{\overline{W}_l \in S_l} \{||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l J_l \overline{W}_l B_{l-1}||_\infty\} \\
\geq\ & \max_{\overline{W}_l \in S_l^1} \{||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l J_l \overline{W}_l B_{l-1}||_\infty\} \\
\geq\ & ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1||_{-\infty} \max_{\overline{W}_l \in S_l} ||J_l \overline{W}_l B_{l-1}||_\infty \\
\geq\ & \gamma \sin(r)cb.
\end{aligned}
$$

Similarly,

$$
\begin{aligned}
& \max_{\overline{W}_{L+1} \in S_{L+1}} \{||M_{L+1}(\overline{W}_{L+1})||_2\} \\
=\ & \max_{\overline{W}_{L+1} \in S_{L+1}} \{||(\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)})B_L||_2\} \\
\geq\ & \max_{\overline{W}_{L+1} \in S_{L+1}} \{||(\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)})B_L||_\infty\} \\
\geq\ & 2\sin(r)b\gamma.
\end{aligned}
$$

Then we obtain the desired lower bound:

$$
\begin{aligned}
& \max_{\widetilde{\mathcal{F}} \in \mathbb{F}} \{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||_2^2\} \\
\geq\ & \max_{\overline{W}_j \in S_j, j \in [L+1]} \{||\nabla(\mathcal{F}_{l_x}(x_0)) - \nabla(\mathcal{F}_{l_2}(x_0))||_2^2 + \sum_{l=1}^{L+1} ||M_l(\overline{W}_l)||_2^2\} \\
\geq\ & ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2 + \gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2).
\end{aligned}
$$

By condition $C_1$ and the lower bound just obtained, we have

$$
\begin{aligned}
& \min_{\widetilde{\mathcal{F}} \in \mathbb{F}} \{\frac{||\nabla(\mathcal{F}_{l_x}(x_0)) - \nabla(\mathcal{F}_{l_2}(x_0))||_2^2}{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x_0)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x_0))||_2^2}\} \\
\leq\ & \frac{||\nabla(\mathcal{F}_{l_x}(x_0)) - \nabla(\mathcal{F}_{l_2}(x_0))||_2^2}{||\nabla(\mathcal{F}_{l_x}(x_0)) - \nabla(\mathcal{F}_{l_2}(x_0))||_2^2 + \gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)} \\
=\ & 1 - \frac{\gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)}{||\nabla(\mathcal{F}_{l_x}(x_0)) - \nabla(\mathcal{F}_{l_2}(x_0))||_2^2 + \gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)} \\
\leq\ & 1 - \frac{\gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)}{4A + \gamma^2((L-1)(\sin(r)cb)^2 + c^2 + (2\sin(r)b)^2)}.
\end{aligned}
$$

The theorem is proved. $\square$

We now prove Theorem 4.4.

*Proof.* The proof is similar to that of Theorem 4.3, so certain details are omitted. Let $T_l = \{\mathcal{F}^{l-1}(x) \,|\, x \in S\} \subset \mathbb{R}$, and $S_l$ the set of $Q \in \mathbb{R}^{n \times n}$ such that $||Q||_\infty \leq \gamma$ and $Qt = 0$ for all $t \in T_l$ and $l \in [L+1]$. $S_l$ must contain non-zero elements because of condition $C_4$.

Let $\mathbb{F}$ be the set of networks $\widetilde{\mathcal{F}} \in \mathbb{R}^n \to \mathbb{R}^m$

$$\widetilde{\mathcal{F}}(x) = \widetilde{W}_{L+1} \sigma(\widetilde{W}_L \sigma(\widetilde{W}_{L-1} \sigma(\ldots \sigma(\widetilde{W}_1 x))))$$

where $\widetilde{W}_l - W_l \in S_l$. Then for all $\widetilde{\mathcal{F}} \in \mathbb{F}$, we have $\widetilde{\mathcal{F}}^l(x) = \mathcal{F}^l(x)$ for $l \in [L+1]$ and $x \in S$, so $\mathbb{F} \subset H(\gamma)$. As a consequence,

$$
\begin{aligned}
& \int_{x \sim D_x} \min_{l \neq l_x} \{||\mathcal{F}_{l_x}(x) - \mathcal{F}_l(x)||_2^2 I(\mathcal{F}_{l_x}(x) > \mathcal{F}_l(x))\} \mathrm{d}x \\
=\ & \int_{x \sim D_x} \min_{l \neq l_x} \{||\widetilde{\mathcal{F}}_{l_x}(x) - \widetilde{\mathcal{F}}_l(x)||_2^2 I(\widetilde{\mathcal{F}}_{l_x}(x) > \widetilde{\mathcal{F}}_l(x))\} \mathrm{d}x.
\end{aligned}
$$

Let $l_2 = \arg\max_{l \neq l_x}\{||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_l(x))||_2^2\}$. Then

$$
\begin{aligned}
&\int_{x \sim D_x} \max_{l \neq l_x}\{||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_i(x))||_2^2\}\mathrm{d}x \\
=\ &\int_{x \sim D_x} ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2\mathrm{d}x
\end{aligned}
$$

and

$$
\begin{aligned}
&\int_{x \sim D_x} \max_{l \neq l_x}\{||\nabla(\widetilde{\mathcal{F}}_{l_x}(x)) - \nabla(\widetilde{\mathcal{F}}_i(x))||_2^2\}\mathrm{d}x \\
\geq\ &\int_{x \sim D_x} ||\nabla(\widetilde{\mathcal{F}}_{l_x}(x)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x))||_2^2\mathrm{d}x.
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
&\min_{\widetilde{\mathcal{F}} \in H(\gamma)}\{\tfrac{R(\widetilde{\mathcal{F}}, D_x)}{R(\mathcal{F}, D_x)}\} \\
\leq\ &\min_{\widetilde{\mathcal{F}} \in \mathbb{F}}\{\tfrac{R(\widetilde{\mathcal{F}}, D_x)}{R(\mathcal{F}, D_x)}\} \\
\leq\ &\min_{\widetilde{\mathcal{F}} \in \mathbb{F}}\{\tfrac{\int_{x \sim D_x} ||\mathcal{F}_{l_x}(x) - \mathcal{F}_{l_2}(x)||_2^2\mathrm{d}x}{\int_{x \sim D_x} ||\widetilde{\mathcal{F}}_{l_x}(x) - \widetilde{\mathcal{F}}_{l_2}(x)||_2^2\mathrm{d}x}\}.
\end{aligned}
$$

We will estimate $\max_{\widetilde{\mathcal{F}} \in \mathbb{F}}\{\int_{x \sim D_x} ||\widetilde{\mathcal{F}}_{l_x}(x) - \widetilde{\mathcal{F}}_{l_2}(x)||_2^2\mathrm{d}x\}$. Let $J_l(x) = \mathrm{diag}(\mathrm{sign}(\mathcal{F}^l(x))) \in \mathbb{R}^{n \times n}$, where $l \in [L]$. Then $\frac{\nabla \mathcal{F}_i(x)}{\nabla x} = W_{L+1}^i(J_L(x)W_L)(J_{L-1}(x)W_{L-1}) \ldots (J_1(x)W_1)$. Also, for all $\widetilde{\mathcal{F}} \in \mathbb{F}$,

$$
\frac{\nabla \widetilde{\mathcal{F}}_i(x)}{\nabla x} = \widetilde{W}_{L+1}^i(J_L(x)\widetilde{W}_L)(J_{L-1}(x)\widetilde{W}_{L-1}) \ldots (J_1(x)\widetilde{W}_1).
$$

Denote $\overline{W}_i = \widetilde{W}_i - W_i \in S_i$. Then

$$
\begin{aligned}
&\nabla(\widetilde{\mathcal{F}}_{l_x}(x)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x)) \\
=\ &(\widetilde{W}_{L+1}^{(l_x)} - \widetilde{W}_{L+1}^{(l_2)})(J_L(x)\widetilde{W}_L)(J_{L-1}(x)\widetilde{W}_{L-1}) \ldots (J_1(x)\widetilde{W}_1) \\
=\ &(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)} + \overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)})(J_L(x)(W_L + \overline{W}_L))(J_{L-1}(x)(W_{L-1} + \overline{W}_{L-1})) \ldots (J_1(x)(W_1 + \overline{W}_1))
\end{aligned}
$$

Let $A_l(x) = \frac{\nabla \mathcal{F}^L(t)}{\nabla \mathcal{F}^l(t)}|_{t=x}$, $B_l(x) = \frac{\nabla \mathcal{F}^l(t)}{\nabla t}|_{t=x}$ where $l \in [L]$. Then

$A_l(x) = (J_L(x)W_L)(J_{L-1}(x)W_{L-1}) \ldots (J_{l+1}(x)W_{l+1})$ and

$B_l(x) = (J_l(x)W_i)(J_{l-1}(x)W_{l-1}) \ldots (J_1(x)W_1).$

Let $\gamma(x, \overline{W}_1) = (W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)J_1(x)\overline{W}_1$, $M_l(x, \overline{W}_l) = (W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l(x)J_l(x)\overline{W}_l B_{l-1}(x)$ where $l \in \{2, 3, \ldots, L\}$, $M_{L+1}(x, \overline{W}_{L+1}) = (\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)})B_L(x)$. By Lemma 6.3,

$$
\begin{aligned}
&\max_{\widetilde{\mathcal{F}} \in \mathbb{F}}\{\int_{x \sim D_x} ||\nabla(\widetilde{\mathcal{F}}_{l_x}(x)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x))||_2^2\mathrm{d}x\} \\
\geq\ &\max_{\overline{W}_l \in S_l, l \in [L+1]}\{\int_{x \sim D_x} ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2 + \sum_{l=1}^{L+1}||M_l(\overline{W}_l)||_2^2\mathrm{d}x\} \\
=\ &\int_{x \sim D_x} ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2 + \sum_{l=1}^{L+1}\max_{\overline{W}_l \in S_l}\{\int_{x \sim D_x} ||M_l(\overline{W}_l)||_2^2\mathrm{d}x\}.
\end{aligned}
$$

Let $\overline{W}_1(k) \in \mathbb{R}^{n \times n}$ be the matrix whose $k$-th row is equal to $k$-th row of $\overline{W}_1$, and other rows are 0. Let $(J_1(x))^{(k)}$ be the $k$-th row of $J_1(x)$. Then

$$
\begin{aligned}
&\max_{\overline{W}_1 \in S_1}\{\int_{x \sim D_x} ||\gamma(x, \overline{W}_1)||_2^2\} \\
=\ &\max_{\overline{W}_1 \in S_1}\{\int_{x \sim D_x} ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)J_1(x)\overline{W}_1||_2^2\mathrm{d}x\} \\
\geq\ &\max_{\overline{W}_1 \in S_1}\{\int_{x \sim D_x} ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)J_1(x)\overline{W}_1||_\infty^2\mathrm{d}x\} \\
\geq\ &\max_{\overline{W}_1 \in S_1, k \in [n]}\{\int_{x \sim D_x} (||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)||_{-\infty}||J_1(x)\overline{W}_1(k)||_\infty)^2\mathrm{d}x\} \\
\geq\ &\max_{k \in [n]}\{\int_{x \sim D_x} (||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)||_{-\infty}I((J_1(x))^{(k)} \neq 0)\gamma)^2\mathrm{d}x\}.
\end{aligned}
$$

By condition $C_2$, we know that $P_{x \sim D_x}(||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)||_{-\infty} > c_1) > \alpha_1$, and by condition $C_4$ and the principle of drawer, there exists a $k \in [n]$ such that $P(x \sim D_x)(I((J_1(x))^{(k)} \neq 0) \mid ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)||_{-\infty} > c) > \gamma$. Thus there exists a $k \in [n]$ such that

$$P(x \sim D_x)(I((J_1(x))^{(k)} \neq 0), ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)||_{-\infty} > c) > \gamma\alpha_1.$$

Then

$$\max_{\overline{W}_1 \in S_1}\{\int_{x \sim D_x} ||\gamma(x, \overline{W}_1)||_2^2\}$$
$$\geq \max_{k \in [n]}\{\int_{x \sim D_x}(||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)||_{-\infty}I((J_1(x))^{(k)} \neq 0)\gamma)^2 dx\}$$
$$\geq \max_{k \in [n]}\{P(x \sim D_x)(I((J_1(x))^{(k)} \neq 0), ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_1(x)||_{-\infty} > c)(c\gamma)^2\}$$
$$\geq (c_1\gamma)^2\gamma\alpha_1.$$

Let $\widetilde{S}_l = \{x \in S_l \mid \text{only one row of } x \text{ is not } 0\}$. Then for $l \in \{2, 3, \ldots, L\}$, we have

$$\max_{\overline{W}_l \in S_l}\{\int_{x \sim D_x} ||M_l(x, \overline{W}_l)||_2^2 dx\}$$
$$= \max_{\overline{W}_l \in S_l}\{\int_{x \sim D_x} ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l(x)J_l(x)\overline{W}_l B_{l-1}(x)||_2^2 dx\}$$
$$\geq \max_{\overline{W}_l \in \widetilde{S}_l}\{\int_{x \sim D_x}(||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l(x)||_{-\infty}||J_l(x)\overline{W}_l B_{l-1}(x)||_\infty)^2 dx\}$$
$$\geq \max_{\overline{W}_l \in \widetilde{S}_l}\{P_{x \sim D_x}(||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l(x)||_{-\infty} > c_l, \ ||J_l(x)\overline{W}_l B_{l-1}(x)||_\infty \geq d_l||J_l(x)\overline{W}_l||_\infty,$$
$$J_l(x)\overline{W}_l \neq 0)(\gamma c_l d_{l-1})^2\}$$

By conditions $C_3$ and $C_2$, we have $P_{x \sim D_x}(||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l(x)||_{-\infty} > c_l, \ ||J_l(x)\overline{W}_l B_{l-1}(x)||_\infty \geq d_l||J_l(x)\overline{W}_l||_\infty) > \alpha_l + \beta_{l-1} - 1$. By condition $C_4$ and the principle of drawer, there exists a $\overline{W}_i \in \widetilde{S}_i$ such that $P_{x \sim D_x}(J_l(x)\overline{W}_l \neq 0 \mid ||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l(x)||_{-\infty} > c_l, \ ||J_l(x)\overline{W}_l B_{l-1}(x)||_\infty \geq d_l||J_l(x)\overline{W}_i||_\infty) > \gamma$. So, there exists a $\overline{W}_l \in \widetilde{S}_l$ such that

$$P_{x \sim D_x}(||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_l(x)||_{-\infty} > c_l, ||J_l(x)\overline{W}_l B_{l-1}(x)||_\infty \geq d_l||J_l(x)\overline{W}_l||_\infty, J_l(x)\overline{W}_l \neq 0)$$
$$> \gamma(\alpha_l + \beta_{l-1} - 1).$$

Then

$$\max_{\overline{W}_i \in S_i}\{\int_{x \sim D_x} ||M_i(x, \overline{W}_i)||_2^2 dx\}$$
$$\geq \max_{\overline{W}_i \in \widetilde{S}_i}\{P_{x \sim D_x}(||(W_{L+1}^{(l_x)} - W_{L+1}^{(l_2)})A_i(x)||_{-\infty} > c_i,$$
$$||J_i(x)\overline{W}_i B_{i-1}(x)||_\infty \geq d_i||J_i(x)\overline{W}_i||_\infty, \ J_i(x)\overline{W}_i \neq 0)(\gamma c_i d_{i-1})^2\}$$
$$\geq (\gamma c_i d_{i-1})^2\gamma(\alpha_i + \beta_{i-1} - 1).$$

Similarly, by condition $C_3$, we have

$$\max_{\overline{W}_{L+1} \in S_{L+1}}\{\int_{x \sim D_x} ||M_{L+1}(x, \overline{W}_{L+1})||_2^2 dx\}$$
$$= \max_{\overline{W}_{L+1} \in S_{L+1}}\{\int_{x \sim D_x} ||(\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)})B_L(x)||_2^2 dx\}$$
$$\geq \max_{\overline{W}_{L+1} \in S_{L+1}}\{\int_{x \sim D_x} ||(\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)})B_L(x)||_\infty^2 dx\}$$
$$\geq \max_{\overline{W}_{L+1} \in S_{L+1}}\{\int_{x \sim D_x} I(||(\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)})B_L(x)||_\infty \geq d_L||\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)}||_\infty)(d_L||\overline{W}_{L+1}^{(l_x)} - \overline{W}_{L+1}^{(l_2)}||_\infty)^2 dx\}$$
$$\geq \beta_L(d_L\gamma)^2.$$

25

Then we have the desired lower bound

$$\max_{\widetilde{\mathcal{F}} \in \mathbb{F}} \{ \int_{x \sim D_x} ||\nabla(\widetilde{\mathcal{F}}_{l_x}(x)) - \nabla(\widetilde{\mathcal{F}}_{l_2}(x))||_2^2 dx \}$$

$$\geq \int_{x \sim D_x} ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2 dx + \sum_{i=1}^{L+1} \max_{\overline{W}_i \in S_i} \{ \int_{x \sim D_x} ||M_i(\overline{W}_i)||_2^2 dx \}$$

$$\geq \int_{x \sim D_x} ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2 dx + (\gamma c)^2 \alpha_1 \gamma + \sum_{i=2}^{L} (\gamma c_i d_{i-1})^2 \gamma(\alpha_i + \beta_{i-1} - 1) + \beta_L (d_L \gamma)^2.$$

By condition $C_1$ and the lower bound just obtained, we have

$$\min_{\widetilde{\mathcal{F}} \in \mathbb{F}} \left\{ \frac{\int_{x \sim D_x} ||\mathcal{F}_{l_x}(x) - \mathcal{F}_{l_2}(x)||_2^2 dx}{\int_{x \sim D_x} ||\widetilde{\mathcal{F}}_{l_x}(x) - \widetilde{\mathcal{F}}_{l_2}(x)||_2^2 dx} \right\}$$

$$\leq \frac{\int_{x \sim D_x} ||\mathcal{F}_{l_x}(x) - \mathcal{F}_{l_2}(x)||_2^2 dx}{\int_{x \sim D_x} ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2 dx + (\gamma c_1)^2 \alpha_1 \gamma_1 + \sum_{i=2}^{L} (\gamma c_i d_{i-1})^2 \gamma_i (\alpha_i + \beta_{i-1} - 1) + \beta_L (d_L \gamma)^2}$$

$$= 1 - \frac{(\gamma c_1)^2 \alpha_1 \gamma_1 + \sum_{i=2}^{L} (\gamma c_i d_{i-1})^2 \gamma_i (\alpha_i + \beta_{i-1} - 1) + \beta_L (d_L \gamma)^2}{\int_{x \sim D_x} ||\nabla(\mathcal{F}_{l_x}(x)) - \nabla(\mathcal{F}_{l_2}(x))||_2^2 dx + (\gamma c_1)^2 \alpha_1 \gamma_1 + \sum_{i=2}^{L} (\gamma c_i d_{i-1})^2 \gamma_i (\alpha_i + \beta_{i-1} - 1) + \beta_L (d_L \gamma)^2}$$

$$\leq 1 - \frac{(\gamma c_1)^2 \alpha_1 \gamma_1 + \sum_{i=2}^{L} (\gamma c_i d_{i-1})^2 \gamma_i (\alpha_i + \beta_{i-1} - 1) + \beta_L (d_L \gamma)^2}{4A + (\gamma c_1)^2 \alpha_1 \gamma_1 + \sum_{i=2}^{L} (\gamma c_i d_{i-1})^2 \gamma_i (\alpha_i + \beta_{i-1} - 1) + \beta_L (d_L \gamma)^2}.$$

The theorem is proved. $\square$

# 7  Conclusion

The adversarial parameter attack for DNNs is proposed. In the attack, the adversary makes small changes to the parameters of a trained DNN such that the attacked DNN will keep the accuracy of the original DNN as much as possible, but makes the robustness as low as possible. The goal of the attack is that the attacked DNN is imperceptible to the user and at the same time the robustness of the DNN is broken. The existence of adversarial parameters is proved under certain conditions and effective adversarial parameter attack algorithms are also given.

In general, it is still out of reach to provide provable safety DNNs in real-world applications, and one of the ways to develop safer DNN models and training methods, and evaluate the safety of the trained model against existing attack methods. In other words, a DNN to be deployed is considered safe if it is safe against existing attacks in certain sense. From this viewpoint, it is valuable to have more attack methods. This is similar to the cryptanalysis [13], where much more matured theory and attack methods are developed.

# References

[1] N. Akhtar and A. Mian. Threat of Adversarial Attacks on Deep Learning in Computer Vision: A Survey. arXiv:1801.00553v3, 2018.

[2] B. Allen, S. Agarwal, J. Kalpathy-Cramer, K. Dreyer. Democratizing AI. *J. Am. Coll. Radiol.*, 16(7), 961-963, 2019.

[3] A. Azulay and Y. Weiss. Why Do Deep Convolutional Networks Generalize so Poorly to Small Image Transformations? *Journal of Machine Learning Research*, 20, 1-25, 2019.

[4] A. Athalye, N. Carlini, D. Wagner. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. *Proc. ICML*, PMLR, 2018: 274-283.

[5] T. Bai, J. Luo, J. Zhao. Recent Advances in Understanding Adversarial Robustness of Deep Neural Networks. arXiv:2011.01539, 2020.

[6] A. Bastounis, A.C. Hansen, V. Vlačić. The Mathematics of Adversarial Attacks in AI - Why Deep Learning is Unstable Despite the Existence of Stable Neural Networks. arXiv:2109.06098, 2021.

[7] J. Breier, X. Hou, D. Jap, L. Ma, S. Bhasin, Y. Liu. DeepLaser: Practical Fault Attack on Deep Neural Networks. arXiv:1806.05859, 2018.

[8] J. Cohen, E. Rosenfeld, Z. Kolter. Certified Adversarial Robustness via Randomized Smoothing. *Proc. ICML'2019*, PMLR, 1310-1320, 2019.

[9] G. Cybenko. Approximation by Superpositions of a Sigmoidal Function. *Mathematics of Control, Signals and Systems*, 2(4): 303-314, 1989.

[10] C. Etmann, S. Lunz, P. Maass, C.B. Schonlieb. On the Connection Between Adversarial Robustness and Saliency Map Interpretability. arXiv:1905.04172, 2019.

[11] K. He, X. Zhang, S. Ren, J. Sun. Deep Residual Learning for Image Recognition. *Proc. CVPR*, 770-778, 2016.

[12] M. Hein and M. Andriushchenko. Formal Guarantees on the Robustness of a Classifier Against Adversarial Manipulation. *Proc. NIPS*, 2266-2276, 2017.

[13] O. Goldreich. *Foundations of Cryptography, Volume II, Basic Tools*. Cambridge University Press, 2009.

[14] I.J. Goodfellow, J. Shlens, C. Szegedy. Explaining and Harnessing Adversarial Examples. arXiv:1412.6572, 2014.

[15] Y. LeCun, Y. Bengio, G. Hinton. Deep Learning. *Nature*, 521(7553), 436-444, 2015.

[16] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner. Gradient-based Learning Applied to Document Recognition. *Proc. of the IEEE*, 86(11): 2278-2324, 1998.

[17] H. Li, A. Kadav, I. Durdanovic, H. Samet, H.P. Graf. Pruning Filters for Efficient Convnets. arXiv:1608.08710, 2016.

[18] Y. Liu, L. Wei, B. Luo, Q. Xu. Fault Injection Attack on Deep Neural Network. *Proc. of the IEEE/ACM International Conference on Computer-Aided Design*, 131-138, 2017.

[19] A. Ma, F. Faghri, N. Papernot, A.M. Farahmand. SOAR: Second-Order Adversarial Regularization. arXiv:2004.01832, 2020.

[20] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu. Towards Deep Learning Models Resistant to Adversarial Attacks. arXiv:1706.06083, 2017.

[21] A.S. Morcos, D.G.T. Barrett, NC. Rabinowitz, M. Botvinick. On the Importance of Single Directions for Generalization. t arXiv:1803.06959, 2018.

[22] B. Neyshabur, R. Tomioka, N. Srebro. Norm-based Capacity Control in Neural Networks. *Proc. COLT'15*, 1376-1401, 2015.

[23] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z.B. Celik, A. Swami. The Limitations of Deep Learning in Adversarial Settings. *IEEE European Symposium on Security and Privacy*, IEEE Press, 2016: 372-387.

[24] A. Raghunathan, J. Steinhardt, P. Liang. Certified Defenses Against Adversarial Examples. ArXiv: 1801.09344, 2018.

[25] A. RoyChowdhury, P. Sharma, E. Learned-Miller. Reducing Duplicate Filters in Deep Neural Networks. *NIPS workshop on Deep Learning: Bridging Theory and Practice*, 1, 2017.

[26] W. Shang, K. Sohn, D. Almeida, H. Lee. Understanding and Improving Convolutional Neural Networks via Concatenated Rectified Linear Units. *Proc. ICML*, PMLR, 2217-2225, 2016.

[27] A. Shafahi, W.R. Huang, C. Studer, S. Feizi, T. Goldstein. Are Adversarial Examples Inevitable? arXiv:1809.02104, 2018.

[28] C.J. Simon-Gabriel, Y. Ollivier, L. Bottou, B. Schölkopf, D. Lopez-Paz. First-order Adversarial Vulnerability of Neural Networks and Input Dimension. *Proc. ICML*, PMLR, 5809-5817, 2019.

[29] K. Simonyan and A. Zisserman. Very Deep Convolutional Networks for Large-scale Image Recognition. arXiv:1409.1556, 2014.

[30] X. Sun, Z. Zhang, X. Ren, R. Luo, L. Li. Exploring the Vulnerability of Deep Neural Networks: A Study of Parameter Corruption. arXiv:2006.05620, 2020.

[31] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I.J. Goodfellow, R. Fergus. Intriguing Properties of Neural Networks. arXiv:1312.6199, 2013.

[32] F. Tramèr, N. Papernot, I. Goodfellow, D. Boneh, P. McDaniel. The Space of Transferable Adversarial Examples. arXiv:1704.03453, 2017.

[33] Y.L. Tsai, C.Y. Hsu, C.M. Yu, P.Y. Chen. Formalizing Generalization and Robustness of Neural Networks to Weight Perturbations. arXiv:2103.02200, 2021.

[34] I.Y. Tyukin, D.J. Higham, A.N. Gorban. On Adversarial Examples and Stealth Attacks in Artificial Intelligence Systems. *2020 International Joint Conference on Neural Networks*, 1-6, IEEE Press, 2020.

[35] I.Y. Tyukin, D.J. Higham, A.N. Gorban. E. Woldegeorgis. The Feasibility and Inevitability of Stealth Attacks. arXiv2106.13997, 2021.

[36] Z. Wang, C. Xiang, W. Zou, C. Xu. MMA Regularization: Decorrelating Weights of Neural Networks by Maximizing the Minimal Angles. arXiv:2006.06527, 2020.

[37] T.W. Weng, P. Zhao, S. Liu, P.Y. Chen, X. Lin, L. Daniel. Towards Certificated Model Robustness Against Weight Perturbations. *Proc. of the AAAI*, 34(04): 6356-6363, 2020.

[38] H. Xu, Y. Ma, H.C. Liu, D, Deb, H. Liu J.L. Tang, A.K. Jain. Adversarial Attacks and Defenses in Images, Graphs and Text: A Review. *International Journal of Automation and Computing*, 17(2), 151-178, 2020.

[39] L. Yu and X.S. Gao. Improve the Robustness and Accuracy of Deep Neural Network with $L_{2,\infty}$ Normalization, arXiv:2010.04912, 2020.

[40] L. Yu and X.S. Gao. Robust and Information-theoretically Safe Bias Classifier against Adversarial Attacks. arXiv:2111.04404, 2021.

[41] H. Zhang, Y. Yu, J. Jiao, E.P. Xing, L.E. Ghaoui, M.I. Jordan. Theoretically Principled Trade-off Between Robustness and Accuracy. *Proc. ICML*, 2019.

[42] X.Y. Zhang, C.L. Liu, C.Y. Suen. Towards Robust Pattern Recognition: A Review. *Proc. of the IEEE*, 108(6), 894-922, 2020.

[43] P. Zhao, S. Wang, C. Gongye, Y. Wang, Y. Fei, X. Lin. Fault Sneaking Attack: a Stealthy Framework for Misleading Deep Neural Networks. *Proc. of the 56th Annual Design Automation Conference*, 2019.