# On the Convergence and Robustness of Adversarial Training

Yisen Wang [* 1]  Xingjun Ma [* 2]  James Bailey [2]  Jinfeng Yi [1]  Bowen Zhou [1]  Quanquan Gu [3]

## Abstract

Improving the robustness of deep neural networks (DNNs) to adversarial examples is an important yet challenging problem for secure deep learning. Across existing defense techniques, adversarial training with Projected Gradient Decent (PGD) is amongst the most effective. Adversarial training solves a min-max optimization problem, with the *inner maximization* generating adversarial examples by maximizing the classification loss, and the *outer minimization* finding model parameters by minimizing the loss on adversarial examples generated from the inner maximization. A criterion that measures how well the inner maximization is solved is therefore crucial for adversarial training. In this paper, we propose such a criterion, namely First-Order Stationary Condition for constrained optimization (FOSC), to quantitatively evaluate the convergence quality of adversarial examples found in the inner maximization. With FOSC, we find that to ensure better robustness, it is essential to use adversarial examples with better convergence quality at the *later stages* of training. Yet at the early stages, high convergence quality adversarial examples are not necessary and may even lead to poor robustness. Based on these observations, we propose a *dynamic* training strategy to gradually increase the convergence quality of the generated adversarial examples, which significantly improves the robustness of adversarial training. Our theoretical and empirical results show the effectiveness of the proposed method.

## 1. Introduction

Although deep neural networks (DNNs) have achieved great success in a number of fields such as computer vision (He et al., 2016) and natural language processing (Devlin et al., 2018), they are vulnerable to adversarial examples crafted by adding small, human imperceptible adversarial perturbations to normal examples (Szegedy et al., 2013; Goodfellow et al., 2015). Such vulnerability of DNNs raises security concerns about their practicability in security-sensitive applications such as face recognition (Kurakin et al., 2016) and autonomous driving (Chen et al., 2015). Defense techniques that can improve DNN robustness to adversarial examples have thus become crucial for secure deep learning.

There exist several defense techniques (*i.e.*, "defense models"), such as input denoising (Guo et al., 2018), gradient regularization (Papernot et al., 2017), and adversarial training (Madry et al., 2018). However, many of these defense models provide either only marginal robustness or have been evaded by new attacks (Athalye et al., 2018). One defense model that demonstrates moderate robustness, and has thus far not been comprehensively attacked, is *adversarial training* (Athalye et al., 2018). Given a $C$-class dataset $S = \{(\mathbf{x}_i^0, y_i)\}_{i=1}^n$ with $\mathbf{x}_i^0 \in \mathbb{R}^d$ as a normal example in the $d$-dimensional input space and $y_i \in \{1, \cdots, C\}$ as its associated label, the objective of adversarial training is to solve the following *min-max optimization* problem:

$$\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}_i - \mathbf{x}_i^0\|_\infty \leq \epsilon} \ell(h_{\boldsymbol{\theta}}(\mathbf{x}_i), y_i), \qquad (1)$$

where $h_{\boldsymbol{\theta}} : \mathbb{R}^d \to \mathbb{R}^C$ is the DNN function, $\mathbf{x}_i$ is the adversarial example of $\mathbf{x}_i^0$, $\ell(h_{\boldsymbol{\theta}}(\mathbf{x}_i), y_i)$ is the loss function on the adversarial example $(\mathbf{x}_i, y_i)$, and $\epsilon$ is the maximum perturbation constraint[1]. The *inner maximization* problem is to find an adversarial example $\mathbf{x}_i$ within the $\epsilon$-ball around a given normal example $\mathbf{x}_i^0$ (*i.e.*, $\|\mathbf{x}_i - \mathbf{x}_i^0\|_\infty \leq \epsilon$) that maximizes the classification loss $\ell$. It is typically nonconcave with respect to the adversarial example. On the other hand, the *outer minimization* problem is to find model parameters that minimize the loss $\ell$ on adversarial examples $\{\mathbf{x}_i\}_{i=1}^n$ that generated from the inner maximization. This is the problem of training a robust classifier on adversarial examples. Therefore, how well the inner maximization problem is solved directly affects the performance of the outer minimization, *i.e.*, the robustness of the classifier.

---

[*]Equal contribution (Yisen Wang is now at Peking University and Xingjun Ma is now at Deakin University) [1]JD.com [2]The University of Melbourne [3]The University of California, Los Angeles. Correspondence to: Quanquan Gu <qgu@cs.ucla.edu>.

---

[1]We only focus on the infinity norm constraint in this paper, but our algorithms and theory apply to other norms as well.

Several attack methods have been used to solve the inner maximization problem, such as Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015) and Projected Gradient Descent (PGD) (Madry et al., 2018). However, the degree to which they solve the inner maximization problem has not been thoroughly studied. Without an appropriate criterion to measure how well the inner maximization is solved, the adversarial training procedure is difficult to monitor or improve. In this paper, we propose such a criterion, namely First-Order Stationary Condition for constrained optimization (FOSC), to measure the convergence quality of the adversarial examples found in the inner maximization. Our proposed FOSC facilitates monitoring and understanding adversarial training through the lens of convergence quality of the inner maximization, and this in turn motivates us to propose an improved training strategy for better robustness. Our main contributions are as follows:

- We propose a principled criterion FOSC to measure the convergence quality of adversarial examples found in the inner maximization problem of adversarial training. It is well-correlated with the adversarial strength of adversarial examples, and is also a good indicator of the robustness of adversarial training.

- With FOSC, we find that better robustness of adversarial training is associated with training on adversarial examples with better convergence quality in the *later stages*. However, in the early stages, high convergence quality adversarial examples are not necessary and can even be harmful.

- We propose a *dynamic* training strategy to gradually increase the convergence quality of the generated adversarial examples and provide a theoretical guarantee on the overall (min-max) convergence. Experiments show that *dynamic* strategy significantly improves the robustness of adversarial training.

## 2. Related Work

### 2.1. Adversarial Attack

Given a normal example $(\mathbf{x}_i^0, y_i)$ and a DNN $h_{\boldsymbol{\theta}}$, the goal of an attacking method is to find an adversarial example $\mathbf{x}_i$ that remains in the $\epsilon$-ball centered at $\mathbf{x}^0$ ($\|\mathbf{x}_i - \mathbf{x}_i^0\|_\infty \le \epsilon$) but can fool the DNN to make an incorrect prediction ($h_{\boldsymbol{\theta}}(\mathbf{x}_i) \ne y_i$). A wide range of attacking methods have been proposed for the crafting of adversarial examples. Here, we only mention a selection.

**Fast Gradient Sign Method (FGSM).** FGSM perturbs normal examples $\mathbf{x}^0$ for one step ($\mathbf{x}^1$) by the amount $\epsilon$ along the gradient direction (Goodfellow et al., 2015):

$$\mathbf{x}^1 = \mathbf{x}^0 + \epsilon \cdot \mathrm{sign}(\nabla_{\mathbf{x}}\ell(h_{\boldsymbol{\theta}}(\mathbf{x}^0), y)). \qquad (2)$$

**Projected Gradient Descent (PGD).** PGD perturbs normal example $\mathbf{x}^0$ for a number of steps $K$ with smaller step size. After each step of perturbation, PGD projects the adversarial example back onto the $\epsilon$-ball of $\mathbf{x}^0$, if it goes beyond the $\epsilon$-ball (Madry et al., 2018):

$$\mathbf{x}^k = \Pi\big(\mathbf{x}^{k-1} + \alpha \cdot \mathrm{sign}(\nabla_{\mathbf{x}}\ell(h_{\boldsymbol{\theta}}(\mathbf{x}^{k-1}), y))\big), \qquad (3)$$

where $\alpha$ is the step size, $\Pi(\cdot)$ is the projection function, and $\mathbf{x}^k$ is the adversarial example at the $k$-th step.

There are also other types of attacking methods, *e.g.*, Jacobian-based Saliency Map Attack (JSMA) (Papernot et al., 2016a), C&W attack (Carlini & Wagner, 2017) and Frank-Wolfe based attack (Chen et al., 2018). PGD is regarded as the strongest first-order attack, and C&W is among the strongest attacks to date.

### 2.2. Adversarial Defense

A number of defense models have been developed such as defensive distillation (Papernot et al., 2016b), feature analysis (Xu et al., 2017; Ma et al., 2018), input denoising (Guo et al., 2018; Liao et al., 2018; Samangouei et al., 2018), gradient regularization (Gu & Rigazio, 2014; Papernot et al., 2017; Tramèr et al., 2018; Ross & Doshi-Velez, 2018), model compression (Liu et al., 2018; Das et al., 2018; Rakin et al., 2018) and adversarial training (Goodfellow et al., 2015; Nøkland, 2015; Madry et al., 2018), among which adversarial training is the most effective.

Adversarial training improves the model robustness by training on adversarial examples generated by FGSM and PGD (Goodfellow et al., 2015; Madry et al., 2018). Tramèr et al. (2018) proposed an ensemble adversarial training on adversarial examples generated from a number of pretrained models. Kolter & Wong (2018) developed a provable robust model that minimizes worst-case loss over a convex outer region. In a recent study by (Athalye et al., 2018), adversarial training on PGD adversarial examples was demonstrated to be the state-of-of-art defense model. Several improvements of PGD adversarial training have also been proposed, such as Lipschitz regularization (Cisse et al., 2017; Hein & Andriushchenko, 2017; Yan et al., 2018; Farnia et al., 2019), and curriculum adversarial training (Cai et al., 2018).

Despite these studies, a deeper understanding of adversarial training and a clear direction for further improvements is largely missing. The inner maximization problem in Eq. (1) lacks an effective criterion that can quantitatively measure the convergence quality of training adversarial examples generated by different attacking methods (which in turn influences the analysis of the whole min-max problem). In this paper, we propose such a criterion and provide new understanding of the robustness of adversarial training. We design a dynamic training strategy that significantly improves the robustness of the standard PGD adversarial training.

# 3. Evaluation of the Inner Maximization

## 3.1. Quantitative Criterion: FOSC

In Eq. (1), the inner maximization problem is a constrained optimization problem, and is in general globally nonconcave. Since the gradient norm of $h$ is not an appropriate criterion for nonconvex/nonconcave constrained optimization problems, inspired by Frank-Wolfe gap (Frank & Wolfe, 1956), we propose a First-Order Stationary Condition for constrained optimization (FOSC) as the convergence criterion for the inner maximization problem, which is affine invariant and not tied to any specific choice of norm:

$$c(\mathbf{x}^k) = \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x} - \mathbf{x}^k, \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k) \rangle, \qquad (4)$$

where $\mathcal{X} = \{\mathbf{x} | \|\mathbf{x} - \mathbf{x}^0\|_\infty \leq \epsilon\}$ is the input domain of the $\epsilon$-ball around normal example $\mathbf{x}^0$, $f(\boldsymbol{\theta}, \mathbf{x}^k) = \ell(h_{\boldsymbol{\theta}}(\mathbf{x}^k), y)$ and $\langle \cdot \rangle$ is the inner product. Note that $c(\mathbf{x}^k) \geq 0$, and a smaller value of $c(\mathbf{x}^k)$ indicates a better solution of the inner maximization (or equivalently, better convergence quality of the adversarial example $\mathbf{x}^k$).

The criterion FOSC in Eq. (4) can be shown to have the following closed-form solution:

$$\begin{aligned}
c(\mathbf{x}^k) &= \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x} - \mathbf{x}^k, \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k) \rangle \\
&= \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x} - \mathbf{x}^0 + \mathbf{x}^0 - \mathbf{x}^k, \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k) \rangle \\
&= \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x} - \mathbf{x}^0, \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k) \rangle \\
&\quad + \langle \mathbf{x}^k - \mathbf{x}^0, -\nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k) \rangle \\
&= \epsilon \|\nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k)\|_1 - \langle \mathbf{x}^k - \mathbf{x}^0, \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}^k) \rangle.
\end{aligned}$$

As an example-wise criterion, $c(\mathbf{x}^k)$ measures the convergence quality of adversarial example $\mathbf{x}^k$ with respect to both the perturbation constraint and the loss function. Optimal convergence where $c(\mathbf{x}^k) = 0$ can be achieved when 1) $\nabla f(\boldsymbol{\theta}, \mathbf{x}^k) = 0$, i.e., $\mathbf{x}^k$ is a stationary point in the interior of $\mathcal{X}$; or 2) $\mathbf{x}^k - \mathbf{x}^0 = \epsilon \cdot \text{sign}(\nabla f(\boldsymbol{\theta}, \mathbf{x}^k))$, that is, local maximum point of $f(\boldsymbol{\theta}, \mathbf{x}^k)$ is reached on the boundary of $\mathcal{X}$. The proposed criterion FOSC allows the monitoring of convergence quality of the inner maximization problem, and provides a new perspective of adversarial training.

## 3.2. FOSC View of Adversarial Training

In this subsection, we will use FOSC to investigate the robustness and learning process of adversarial training. First though, we investigate its correlation with the traditional measures of accuracy and loss.

**FOSC View of Adversarial Strength.** We train an 8-layer Convolutional Neural Network (CNN) on CIFAR-10 using 10-step PGD (PGD-10) with step size $\epsilon/4$, maximum perturbation $\epsilon = 8/255$, following the standard setting in Madry



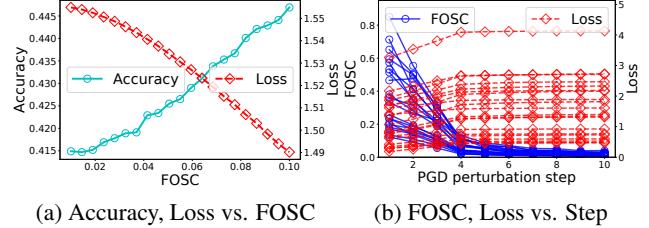(a) Accuracy, Loss vs. FOSC     (b) FOSC, Loss vs. Step

*Figure 1.* The correlation between convergence quality (FOSC) and adversarial strength (accuracy and loss). (a): For PGD-10 CIFAR-10 adversarial examples: the lower the FOSC (x-axis), the lower the accuracy (left y-axis) and the higher the loss (right y-axis). (b): For 20 randomly selected adversarial examples (each line is an example): PGD perturbation step (x-axis) versus FOSC (left y-axis) and loss (right y-axis).

et al. (2018). We then apply the same PGD-10 attack on CIFAR-10 test images to craft adversarial examples, and divide the crafted adversarial examples into 20 consecutive groups of different convergence levels of FOSC value ranging from 0.0 to 0.1. The test accuracy and average loss of adversarial examples in each group are in Figure 1a. We observe FOSC has a linear correlation with both accuracy and loss: the lower the FOSC, the lower (resp. higher) the accuracy (resp. loss).

We further show the intermediate perturbation steps of PGD for 20 randomly selected adversarial examples in Figure 1b. As perturbation step increases, FOSC decreases consistently towards 0, while loss increases and stabilizes at a much wider range of values. Compared to the loss, FOSC provides a *comparable* and *consistent* measurement of adversarial strength: the closer the FOSC to 0, the stronger the attack.

In summary, the proposed FOSC is well correlated with the adversarial strength and also more consistent than the loss, making it a promising tool to monitor adversarial training.

**FOSC View of Adversarial Robustness.** We first investigate the correlation among the final robustness of adversarial training, loss, and FOSC. In particular, we evaluate PGD adversarial training on CIFAR-10 in two settings: 1) varying PGD step size from $\epsilon, \epsilon/2$ to $\epsilon/8$ while fixing step number as 20, and 2) varying PGD step number from 10 to 40 while fixing step size as $\epsilon/6$. In each setting, we cross test (white-box) the robustness of the final model against PGD attacks in the same setting on CIFAR-10 test images. For each defense model, we also compute the distributions of FOSC and loss (using Gaussian kernel density estimation (Parzen, 1962)) for the last epoch generated adversarial examples.

As shown in Figure 2, when varying step size, the best robustness against all test attacks is observed for PGD-$\epsilon/2$ or PGD-$\epsilon/4$ (Figure 2a), of which the FOSC distributions are more concentrated around 0 (Figure 2c) but their loss distributions are almost the same (Figure 2e). When varying step number, the final robustness values are very similar
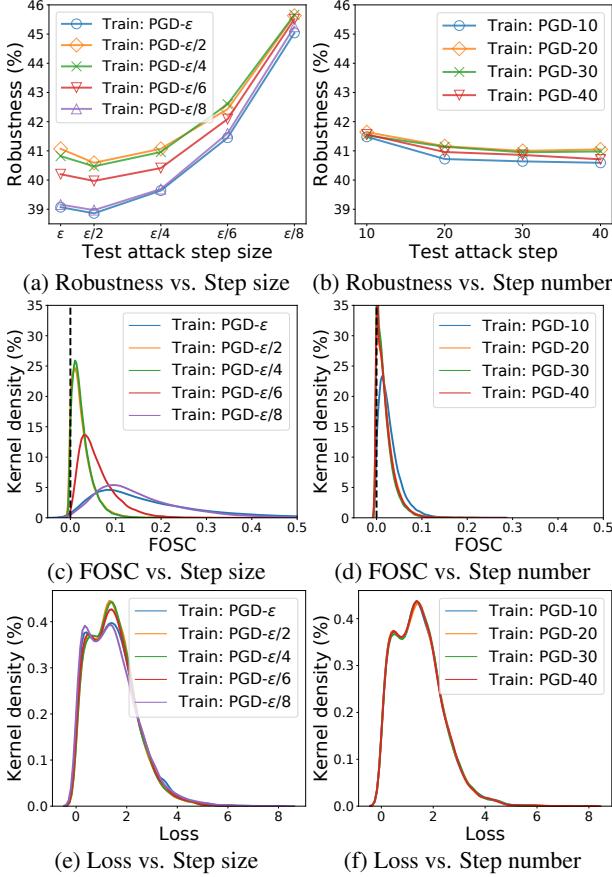
(a) Robustness vs. Step size

(b) Robustness vs. Step number

(c) FOSC vs. Step size

(d) FOSC vs. Step number

(e) Loss vs. Step size

(f) Loss vs. Step number

*Figure 2.* Robustness of PGD adversarial training with (a) varying step size (fixed step number 20), or (b) varying step number (fixed step size $\epsilon/6$). The FOSC distributions (c)/(d) reflect the robustness of adversarial training in (a)/(b), *i.e.*, the lower the FOSC, the better the robustness. The loss distributions (e)/(f) are almost the same for different settings.

(Figure 2b), which is also reflected by the similar FOSC distributions (Figure 2d), but the loss distributions are again almost the same (Figure 2f). Revisiting Figure 2b where the step size is $\epsilon/6$, it is notable that increasing PGD steps only brings marginal or no robustness gain when the steps are more than sufficient to reach the surface of the $\epsilon$-ball: 12 steps of $\epsilon/6$ perturbation following the same gradient direction can reach the surface of the $\epsilon$-ball from any starting point. The above observations indicate that FOSC is a more reliable indicator of the final robustness of PGD adversarial training, compared to the loss.

**Rethinking the Adversarial Training Process.** To provide more insights into the learning process of adversarial training, we show the distributions of FOSC at three distinct learning stages in Figure 3: 1) early stage (epoch 10), 2) middle stage (epoch 60), and 3) later stage (epoch 100) (120 epochs in total). We only focus on two defense models: training with 10-step PGD-$\epsilon/4$ and 10-step PGD-$\epsilon/8$ (the best and worst model observed in Figure 2a respectively).



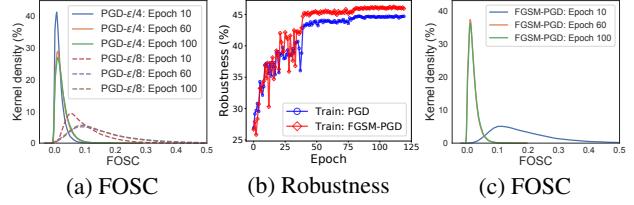(a) FOSC

(b) Robustness

(c) FOSC

*Figure 3.* (a): FOSC distribution at intermediate epochs (10, 60, 100) for adversarial training of 10 steps PGD with step size $\epsilon/4$ (PGD-$\epsilon/4$) and step size $\epsilon/8$ (PGD-$\epsilon/8$); (b): The robustness of training with PGD and training with first FGSM then PGD; (c): FOSC distribution at intermediate epochs for training with first FGSM then PGD. Distributions at the 60-th and 100-th epochs overlap each other for (a)/(c).

In Figure 3a, for both models, FOSC at the early stage is significantly lower than the following two stages. Thus, at the early stage, both models can easily find high convergence quality adversarial examples for training; however, it becomes more difficult to do so at the following stages. This suggests overfitting to strong PGD adversarial examples at the early stage. To verify this, we replace the first 20 epochs of PGD-$\epsilon/4$ training with a much weaker FGSM (1 step perturbation of size $\epsilon$), denoted as "FGSM-PGD", and show its robustness and FOSC distribution in Figure 3b and 3c respectively. We find that by simply using weaker FGSM adversarial examples at the early stage, the final robustness and the convergence quality of adversarial examples found by PGD at the later stage are both significantly improved. The FOSC density between $[0, 0.1]$ is improved to above $35\%$ (green solid line in Figure 3c) from less than $30\%$ (green solid line in Figure 3a). This indicates strong PGD attacks are not necessary for the *early stage* of training, or even deteriorate the robustness. In the next section, we will propose a dynamic training strategy to address this issue.

## 4. Dynamic Adversarial Training

In this section, we first introduce the proposed dynamic adversarial training strategy. Then we provide a theoretical convergence analysis of the min-max problem in Eq. (1) of the proposed approach.

### 4.1. The Proposed Dynamic Training Strategy

As mentioned in Section 3.2, training on adversarial examples of better convergence quality at the *later stages* leads to higher robustness. However, at the *early stages*, training on high convergence quality adversarial examples may not be helpful. Recalling the criterion FOSC proposed in Section 3.1, we have seen that it is strongly correlated with adversarial strength. Thus, it can be used to monitor the strength of adversarial examples at a fine-grained level. Therefore, we propose to train DNNs with adversarial examples of gradually decreasing FOSC value (increasing convergence

quality), so as to ensure that the network is trained on weak adversarial examples at the early stages and strong adversarial examples at the later stages.

Our proposed dynamic adversarial training algorithm is shown in Algorithm 1. The dynamic criterion FOSC $c_t = \max(c_{\max} - t \cdot c_{\max}/T', 0)$ controls the minimum FOSC value (maximum adversarial strength) of the adversarial examples at the $t$-th epoch of training ($T'$ is slightly smaller than total epochs $T$ to ensure the later stage can be trained on criterion 0). In the early stages of training, $c_t$ is close to the maximum FOSC value $c_{\max}$ corresponding to weak adversarial examples, it then decreases linearly towards zero as training progresses[2], and is zero after the $T'$-th epoch of training. We use PGD to generate the training adversarial examples, however, at each perturbation step of PGD, we monitor the FOSC value and stop the perturbation process for adversarial example whose FOSC value is already smaller than $c_t$, enabled by an indicator control vector $V$. $c_{\max}$ can be estimated by the average FOSC value on a batch of weak adversarial examples such as FGSM.

---

**Algorithm 1** Dynamic Adversarial Training

**Input:** Network $h_{\boldsymbol{\theta}}$, training data $S$, initial model parameters $\boldsymbol{\theta}^0$, step size $\eta_t$, mini-batch $\mathcal{B}$, maximum FOSC value $c_{max}$, training epochs $T$, FOSC control epoch $T'$, PGD step $K$, PGD step size $\alpha$, maximum perturbation $\epsilon$.
**for** $t = 0$ **to** $T - 1$ **do**
    $c_t = \max(c_{\max} - t \cdot c_{\max}/T', 0)$
    **for** each batch $\mathbf{x}_{\mathcal{B}}^0$ **do**
        $V = \mathbb{1}_{\mathcal{B}}$     # *control vector of all elements is 1*
        **while** $\sum V > 0$ & $k < K$ **do**
            $\mathbf{x}_{\mathcal{B}}^{k+1} = \mathbf{x}_{\mathcal{B}}^k + V \cdot \alpha \cdot \text{sign}(\nabla_{\mathbf{x}}\ell(h_{\boldsymbol{\theta}}(\mathbf{x}_{\mathcal{B}}^k), y))$
            $\mathbf{x}_{\mathcal{B}}^k = clip(\mathbf{x}_{\mathcal{B}}^k, \mathbf{x}_{\mathcal{B}}^0 - \epsilon, \mathbf{x}_{\mathcal{B}}^0 + \epsilon)$
            $V = \mathbb{1}_{\mathcal{B}}(c(\mathbf{x}_{1\ldots\mathcal{B}}^k) \le c_t)$     # *The element of V becomes 0 at which FOSC is smaller than $c_t$*
        **end while**
        $\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_t \mathbf{g}(\boldsymbol{\theta}^t)$   # $\mathbf{g}(\boldsymbol{\theta}^t)$ : *stochastic gradient*
    **end for**
**end for**

---

## 4.2. Convergence Analysis

We provide a convergence analysis of our proposed dynamic adversarial training approach (as opposed to just the inner maximization problem) for solving the overall min-max optimization problem in Eq. (1). Due to the nonlinearities in DNNs such as ReLU (Nair & Hinton, 2010) and max-pooling functions, the exact assumptions of Danskin's theorem (Danskin, 2012) do not hold. Nevertheless, given the criterion FOSC that ensures an approximate maximizer of the inner maximization problem, we can still provide a

theoretical convergence guarantee.

In detail, let $\mathbf{x}_i^*(\boldsymbol{\theta}) = \text{argmax}_{\mathbf{x}_i \in \mathcal{X}_i} f(\boldsymbol{\theta}, \mathbf{x}_i)$ where $f(\boldsymbol{\theta}, \mathbf{x}) = \ell(h_{\boldsymbol{\theta}}(\mathbf{x}), y)$ is a shorthand notation for the classification loss function, $\mathcal{X}_i = \{\mathbf{x}|\|\mathbf{x} - \mathbf{x}_i^0\|_\infty \le \epsilon\}$, and $\bar{f}_i(\boldsymbol{\theta}) = \max_{\mathbf{x}_i \in \mathcal{X}_i} f(\boldsymbol{\theta}, \mathbf{x}_i) = f(\boldsymbol{\theta}, \mathbf{x}_i^*(\boldsymbol{\theta}))$, then $\hat{\mathbf{x}}_i(\boldsymbol{\theta})$ is a $\delta$-approximate solution to $\mathbf{x}_i^*(\boldsymbol{\theta})$, if it satisfies that

$$c(\hat{\mathbf{x}}_i(\boldsymbol{\theta})) = \max_{\mathbf{x} \in \mathcal{X}_i}\langle \mathbf{x} - \hat{\mathbf{x}}_i(\boldsymbol{\theta}), \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \hat{\mathbf{x}}_i(\boldsymbol{\theta}))\rangle \le \delta. \quad (5)$$

In addition, denote the objective function in Eq. (1) by $L_S(\boldsymbol{\theta})$, and its gradient by $\nabla L_S(\boldsymbol{\theta}) = 1/n \sum_{i=1}^n \nabla \bar{f}_i(\boldsymbol{\theta}) = 1/n \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \mathbf{x}_i^*(\boldsymbol{\theta}))$. Let $\mathbf{g}(\boldsymbol{\theta}) = 1/|\mathcal{B}| \sum_{i \in \mathcal{B}} \nabla \bar{f}_i(\boldsymbol{\theta})$ be the stochastic gradient of $L_S(\boldsymbol{\theta})$, where $\mathcal{B}$ is the mini-batch. We have $\mathbb{E}[\mathbf{g}(\boldsymbol{\theta})] = \nabla L_S(\boldsymbol{\theta})$. Let $\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \hat{\mathbf{x}}(\boldsymbol{\theta}))$ be the gradient of $f(\boldsymbol{\theta}, \hat{\mathbf{x}}(\boldsymbol{\theta}))$ with respect to $\boldsymbol{\theta}$, and $\hat{\mathbf{g}}(\boldsymbol{\theta}) = 1/|\mathcal{B}| \sum_{i \in \mathcal{B}} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \hat{\mathbf{x}}_i(\boldsymbol{\theta}))$ be the approximate stochastic gradient of $L_S(\boldsymbol{\theta})$.

Before we provide the convergence analysis, we first lay out a few assumptions that are needed for our analysis.

**Assumption 1.** *The function $f(\boldsymbol{\theta}; \mathbf{x})$ satisfies the gradient Lipschitz conditions as follows*

$$\sup_{\mathbf{x}} \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \mathbf{x}) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}', \mathbf{x})\|_2 \le L_{\theta\theta}\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2$$

$$\sup_{\boldsymbol{\theta}} \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \mathbf{x}) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, \mathbf{x}')\|_2 \le L_{\theta x}\|\mathbf{x} - \mathbf{x}'\|_2$$

$$\sup_{\mathbf{x}} \|\nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}) - \nabla_{\mathbf{x}} f(\boldsymbol{\theta}', \mathbf{x})\|_2 \le L_{x\theta}\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2,$$

*where $L_{\theta\theta}, L_{\theta x}, L_{x\theta}$ are positive constants.*

Assumption 1 was made in Sinha et al. (2018), which requires the loss function is smooth in the first and second arguments. While ReLU (Nair & Hinton, 2010) is non-differentiable, recent studies (Allen-Zhu et al., 2018; Du et al., 2018; Zou et al., 2018; Cao & Gu, 2019) showed that the loss function of overparamterized deep neural networks is semi-smooth. This helps justify Assumption 1.

**Assumption 2.** *$f(\boldsymbol{\theta}, \mathbf{x})$ is locally $\mu$-strongly concave in $\mathcal{X}_i = \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}_i^0\|_\infty \le \epsilon\}$ for all $i \in [n]$, i.e., for any $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}_i$, it holds that*

$$f(\boldsymbol{\theta}, \mathbf{x}_1) \le f(\boldsymbol{\theta}, \mathbf{x}_2) + \langle \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}_2), \mathbf{x}_1 - \mathbf{x}_2\rangle - \frac{\mu}{2}\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2.$$

Assumption 2 can be verified using the relation between robust optimization and distributional robust optimization (refer to Sinha et al. (2018); Lee & Raginsky (2018)).

**Assumption 3.** *The variance of the stochastic gradient $\mathbf{g}(\boldsymbol{\theta})$ is bounded by a constant $\sigma^2 > 0$,*

$$\mathbb{E}[\|\mathbf{g}(\boldsymbol{\theta}) - \nabla L_S(\boldsymbol{\theta})\|_2^2] \le \sigma^2,$$

*where $\nabla L_S(\boldsymbol{\theta})$ is the full gradient.*

Assumption 3 is a common assumption made for the analysis of stochastic gradient based optimization algorithms.

---

[2]This is only a simple strategy that works well in our experiments and other strategies could also work here.

**Theorem 1.** *Suppose Assumptions 1, 2 and 3 hold. Let* $\Delta = L_S(\boldsymbol{\theta}^0) - \min_{\boldsymbol{\theta}} L_S(\boldsymbol{\theta})$. *If the step size of the outer minimization is set to* $\eta_t = \eta = \sqrt{\Delta/(TL\sigma^2)}$ *and* $T \geq \Delta L/\sigma^2$. *Then the output of Algorithm 1 satisfies*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\big[\|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2\big] \leq 4\sigma\sqrt{\frac{L\Delta}{T}} + \frac{2L_{\theta x}^2\delta}{\mu},$$

*where* $L = (L_{\theta x}L_{x\theta}/\mu + L_{\theta\theta})$.

The complete proof can be found in the supplementary material. Theorem 1 suggests that if the inner maximization is solved up to a precision so that the criterion FOSC is less than $\delta$, Algorithm 1 can converge to a first-order stationary point at a sublinear rate up to a precision of $2L_{\theta x}^2\delta/\mu$. In practice, if $\delta$ is sufficiently small such that $2L_{\theta x}^2\delta/\mu$ is small enough, Algorithm 1 can find a robust model $\boldsymbol{\theta}^T$. This supports the validity of Algorithm 1.

### 4.3. Relation to Curriculum Learning

Curriculum learning (Bengio et al., 2009) is a learning paradigm in which a model learns from easy examples first then gradually learns from more and more difficult examples. For training with normal examples, it has been shown to be able to speed up convergence and improve generalization. This methodology has been adopted in many applications to enhance model training (Kumar et al., 2010; Jiang et al., 2015). The main challenge for curriculum learning is to define a proper criterion to determine the difficulty/hardness of training examples, so as to design a learning curriculum (*i.e.*, a sequential ordering) mechanism. Our proposed criterion FOSC in Section 3.1 can serve as such a difficulty measure for training examples, and our proposed dynamic approach can be regarded as one type of curriculum learning.

Curriculum learning was used in adversarial training in Cai et al. (2018), with the perturbation step of PGD as the difficulty measure. Their assumption is that more perturbation steps indicate stronger adversarial examples. However, this is not a reliable assumption from the FOSC view of the inner maximization problem: more steps may overshoot and result in suboptimal adversarial examples. Empirical comparisons with (Cai et al., 2018) will be shown in Sec. 5.

## 5. Experiments

In this section, we evaluate the robustness of our proposed training strategy (*Dynamic*) compared with several state-of-the-art defense models, in both the white-box and black-box settings on benchmark datasets MNIST and CIFAR-10. We also provide analysis and insights on the robustness of different defense models. For all adversarial examples, we adopt the infinity norm ball as the maximum perturbation constraint (Madry et al., 2018).

**Baselines.** The baseline defense models we use include 1) *Unsecured*: unsecured training on normal examples; 2) *Standard*: standard adversarial training with PGD attacks (Madry et al., 2018); 3) *Curriculum*: curriculum adversarial training with PGD attacks of gradually increasing the number of perturbation steps (Cai et al., 2018).

### 5.1. Robustness Evaluation

**Defense Settings.** For MNIST, defense models use a 4-layer CNN: 2 convolutional layers followed by 2 dense layers. Batch normalization (BatchNorm) (Ioffe & Szegedy, 2015) and max-pooling (MaxPool) are applied after each convolutional layer. For CIFAR-10, defense models adopt an 8-layer CNN architecture: 3 convolutional blocks followed by 2 dense layers, with each convolutional block has 2 convolutional layers. BatchNorm is applied after each convolutional layer, and MaxPool is applied after every convolutional block. Defense models for both MNIST and CIFAR-10 are trained using Stochastic Gradient Descent (SGD) with momentum 0.9, weight decay $10^{-4}$ and an initial learning rate of 0.01. The learning rate is divided by 10 at the 20-th and 40-th epoch for MNIST (50 epochs in total), and at the 60-th and 100-th epoch for CIFAR-10 (120 epochs in total). All images are normalized into [0, 1].

Except the *Unsecured* model, other defense models including our proposed *Dynamic* model are all trained under the same PGD adversarial training scheme: 10-step PGD attack with random start (adding an initial random perturbation of $[-\epsilon, \epsilon]$ to the normal examples before the PGD perturbation) and step size $\epsilon/4$. The maximum perturbation is set to $\epsilon = 0.3$ for MNIST, and $\epsilon = 8/255$ for CIFAR-10, which is a standard setting for adversarial defense (Athalye et al., 2018; Madry et al., 2018). For *Dynamic* model, we set $c_{max} = 0.5$ for both MNIST and CIFAR-10, and $T' = 40$ for MNIST and $T' = 100$ for CIFAR-10. Other parameters of the baselines are configured as per their original papers.

**White-box Robustness.** For MNIST and CIFAR-10, the attacks used for white-box setting are generated from the original test set images by attacking the defense models using 4 attacking methods: FGSM, PGD-10 (10-step PGD), PGD-20 (20-step PGD), and C&W$_\infty$ ($L_\infty$ version of C&W optimized by PGD for 30 steps). In the white-box setting, all attacking methods have full access to the defense model parameters and are constrained by the same maximum perturbation $\epsilon$. We report the classification accuracy of a defense model under white-box attacks as its white-box robustness.

The white-box results are reported in Table 1. On both datasets, the *Unsecured* model achieves the best test accuracy on clean (unperturbed) images. However, it is not robust to adversarial examples — accuracy drops to zero on strong attacks like PGD-10/20 or C&W$_\infty$. The proposed *Dynamic* model almost achieves the best robustness among all

*Table 1.* White-box robustness (accuracy (%) on white-box test attacks) of different defense models on MNIST and CIFAR-10 datasets.

| Defense | MNIST | | | | | CIFAR-10 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Clean | FGSM | PGD-10 | PGD-20 | C&W$_\infty$ | Clean | FGSM | PGD-10 | PGD-20 | C&W$_\infty$ |
| *Unsecured* | **99.20** | 14.04 | 0.0 | 0.0 | 0.0 | **89.39** | 2.2 | 0.0 | 0.0 | 0.0 |
| *Standard* | 97.61 | 94.71 | 91.21 | 90.62 | 91.03 | 66.31 | 48.65 | 44.39 | 40.02 | 36.33 |
| *Curriculum* | 98.62 | **95.51** | 91.24 | 90.65 | 91.12 | 72.40 | 50.47 | 45.54 | 40.12 | 35.77 |
| ***Dynamic*** | 97.96 | 95.34 | **91.63** | **91.27** | **91.47** | 72.17 | **52.81** | **48.06** | **42.40** | **37.26** |

*Table 2.* Black-box robustness (accuracy (%) on black-box test attacks) of different defense models on MNIST and CIFAR-10 datasets.

| Defense | MNIST | | | | CIFAR-10 | | | |
|---|---|---|---|---|---|---|---|---|
| | FGSM | PGD-10 | PGD-20 | C&W$_\infty$ | FGSM | PGD-10 | PGD-20 | C&W$_\infty$ |
| *Standard* | 96.12 | 95.73 | 95.73 | 97.20 | 65.65 | 65.80 | 65.60 | 66.12 |
| *Curriculum* | 96.59 | 95.87 | 96.09 | 97.52 | 71.25 | 71.44 | 71.13 | 71.94 |
| ***Dynamic*** | **97.60** | **97.01** | **96.97** | **98.36** | **71.95** | **72.15** | **72.02** | **72.85** |

the defense models. Comparing the robustness on MNIST and CIFAR-10, the improvements are more significant on CIFAR-10. This may because MNIST consisting of only black-white digits is a relatively simple dataset where different defense models all work comparably well. Compared to *Standard* adversarial training, *Dynamic* training with convergence quality controlled adversarial examples improves the robustness to a certain extent, especially on the more challenging CIFAR-10 with natural images. This robustness gain is possibly limited by the capacity of the small model (only an 8-layer CNN). Thus we shortly show a series of experiments on WideResNet (Zagoruyko & Komodakis, 2016) where the power of the *Dynamic* strategy is fully unleashed. In Table 1, we see that *Curriculum* improves robustness against weak attacks like FGSM but is less effective against strong attacks like PGD/C&W$_\infty$.

**Benchmarking the State-of-the-art on WideResNet.** To analyze the full power of our proposed *Dynamic* training strategy and also benchmark the state-of-the-art robustness on CIFAR-10, we conduct experiments on a large capacity network WideResNet (Zagoruyko & Komodakis, 2016) (10 times wider than standard ResNet (He et al., 2016)), using the same settings as Madry et al. (2018). The WideResNet achieves an accuracy of 95.2% on clean test images of CIFAR-10. For comparison, we include *Madry's* WideResNet adversarial training and the *Curriculum* model. White-box robustness against FGSM, PGD-20 and C&W$_\infty$ attacks is shown in Table 3. Our proposed *Dynamic* model demonstrates a significant boost over *Madry's* WideResNet adversarial training on FGSM and PGD-20, while *Curriculum* model only achieves slight gains respectively. For the strongest attack C&W$_\infty$, *Curriculum*'s robustness decreases by 4% compared to *Madry's*, while *Dynamic* model achieves the highest robustness.

**Black-box Robustness.** Black-box test attacks are gener-

*Table 3.* White-box robustness (%) of different defense models on CIFAR-10 dataset using WideResNet setting in Madry's baselines.

| Defense | Clean | FGSM | PGD-20 | C&W$_\infty$ |
|---|---|---|---|---|
| *Madry's* | **87.3** | 56.1 | 45.8 | 46.8 |
| *Curriculum* | 77.43 | 57.17 | 46.06 | 42.28 |
| ***Dynamic*** | 85.03 | **63.53** | **48.70** | **47.27** |

ated on the original test set images by attacking a surrogate model with architecture that is either a copy of the defense model (for MNIST) or a more complex ResNet-50 (He et al., 2016) model (for CIFAR-10). Both surrogate models are trained separately from the defense models on the original training sets using *Standard* adversarial training (10-step PGD attack with a random start and step size $\epsilon/4$). The attacking methods used here are the same as the white-box evaluation: FGSM, PGD-10, PGD-20, and C&W$_\infty$.

The robustness of different defense models against black-box attacks is reported in Table 2. Again, the proposed *Dynamic* achieves higher robustness than the other defense models. *Curriculum* also demonstrates a clear improvement over *Standard* adversarial training. The distinctive robustness boosts of *Dynamic* and *Curriculum* indicate that training with weak adversarial examples at the early stage can improve the final robustness.

Compared with the white-box robustness in Table 1, all defense models achieve higher robustness against black-box attacks, even the CIFAR-10 black-box attacks which are generated based on a much more complex ResNet-50 network (the defense network is only an 8-layer CNN). This implies that black-box attacks are indeed less powerful than white-box attacks, at least for the tested attacks. It is also observed that robustness tends to increase from weak attacks like FGSM to stronger attacks like C&W$_\infty$. This implies that stronger attacks tend to have less transferability, an observation which is consistent with Madry et al. (2018).

## 5.2. Further Analysis

**Different Maximum Perturbation Constraints** $\epsilon$**.** We analyze the robustness of defense models *Standard*, *Curriculum* and the proposed *Dynamic*, under different maximum perturbation constraints $\epsilon$ on CIFAR-10. For efficiency, we use the same 8-layer CNN defense architecture as in Sec. 5.1. We see in Figure 4a the white-box robustness of defense models trained with $\epsilon = 8/255$ against different PGD-10 attacks with varying $\epsilon \in [2/255, 8/255]$. *Curriculum* and *Dynamic* models substantially improve the robustness of *Standard* adversarial training, a result consistent with Sec. 5.1. *Dynamic* training is better against stronger attacks with larger perturbations ($\epsilon = 8/255$) than *Curriculum*. *Curriculum* is effective on attacks with smaller perturbations ($\epsilon = 4/255, 2/255$), as similar performance to *Dynamic*. We also train the defense models with different $\epsilon \in [2/255, 8/255]$, and then test their white-box robustness under the same $\epsilon$ (all defense models will tend to have similar low robustness if testing $\epsilon$ is larger than training $\epsilon$). As illustrated in Figure 4b, training with weak attacks at the early stages might have a limit: the robustness gain tends to decrease when the maximum perturbation decreases to $\epsilon = 2/255$. This is not surprising given the fact that the inner maximization problem of adversarial training becomes more concave and easier to solve given the smaller $\epsilon$-ball. However, robustness for this extremely small scale perturbation is arguably less interesting for secure deep learning.
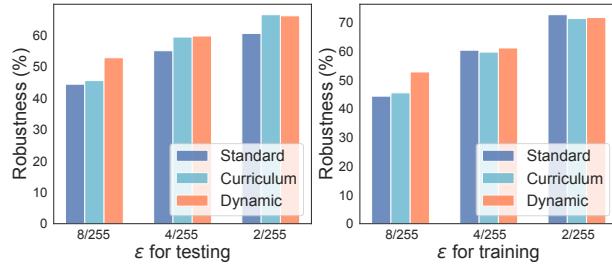


(a) Robustness vs. testing $\epsilon$     (b) Robustness vs. training $\epsilon$

*Figure 4.* (a): White-box robustness on PGD-10 attacks with different testing $\epsilon \in [2/255, 8/255]$; (b): White-box robustness of defense models trained on PGD-10 with different training $\epsilon \in [2/255, 8/255]$.

**Adversarial Training Process.** To understand the learning dynamics of the 3 defense models, we plot the distribution of FOSC at different training epochs in Figures 5. We choose epoch 10/60/100 for *Standard* and *Dynamic*, and epoch 60/90/120 for *Curriculum* as it trains without adversarial examples at early epochs. We see both *Curriculum* and *Dynamic* learn with adversarial examples that are of increasing convergence quality (decreasing FOSC). The difference is that *Dynamic* has more precise control over the convergence quality due to its use of the proposed criterion FOSC, demonstrating more concentrated FOSC distributions which are more separated at different stages of training. While for
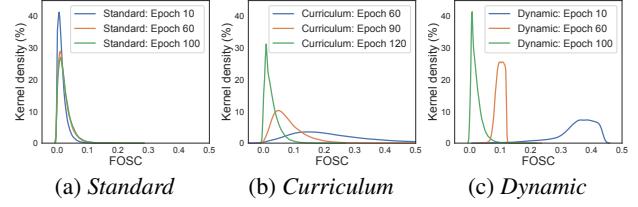


*Figure 5.* The distributions of FOSC at different epochs of training on CIFAR-10 with 10-step PGD of step size $\epsilon/4$ and $\epsilon = 8/255$.

*Curriculum*, the convergence quality of adversarial examples generated by the same number of perturbation steps can span a wide range of values (*e.g.* the flat blue line in Figure 5b), having both weak and strong adversarial examples. Regarding the later stages of training (epoch 100/120), we see *Dynamic* ends up with the best convergence quality (FOSC density over 40% in Figure 5c) followed by *Curriculum* (FOSC density over 30% in Figure 5b) and *Standard* (FOSC density less than 30% in Figure 5a), which is well aligned with their final robustness reported in Tables 1 and 2.

## 6. Discussion and Conclusion

In this paper, we proposed a criterion, First-Order Stationary Condition for constrained optimization (FOSC), to measure the convergence quality of adversarial examples found in the inner maximization of adversarial training. The proposed criterion FOSC is well correlated with adversarial strength and is more consistent than the loss. Based on FOSC, we found that higher robustness of adversarial training can be achieved by training on better convergence quality adversarial examples at the *later stages*, rather than at the *early stages*. Following that, we proposed a *dynamic* training strategy and proved the convergence of the proposed approach for the overall min-max optimization problem under certain assumptions. On benchmark datasets, especially on CIFAR-10 under the WideResNet architecture for attacks with maximum perturbation constraint $\epsilon = 8/255$, our proposed dynamic strategy achieved a significant robustness gain against Madry's state-of-the-art baselines.

Our findings imply that including very hard adversarial examples too early in training possibly inhibits DNN feature learning or encourages premature learning of overly complex features that provide less compression of patterns in the data. Experimental evidences also suggest that the later stages are more correlated with the final robustness, while the early stages are more associated with generalization. Therefore, we conjecture that higher robustness can be obtained by further increasing the diversity of weak adversarial examples in the early stages or generating more powerful adversarial examples in the later stages. The precise characterization of how the early and later stages interact with each other is still an open problem. We believe further exploration of this direction will lead to more robust models.

## Acknowledgements

## References

Allen-Zhu, Z., Li, Y., and Song, Z. A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018.

Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *ICML*, 2018.

Bengio, Y., Louradour, J., Collobert, R., and Weston, J. Curriculum learning. In *ICML*, 2009.

Cai, Q.-Z., Du, M., Liu, C., and Song, D. Curriculum adversarial training. In *IJCAI*, 2018.

Cao, Y. and Gu, Q. A generalization theory of gradient descent for learning over-parameterized deep relu networks. *arXiv preprint arXiv:1902.01384*, 2019.

Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy*, 2017.

Chen, C., Seff, A., Kornhauser, A., and Xiao, J. Deepdriving: Learning affordance for direct perception in autonomous driving. In *CVPR*, 2015.

Chen, J., Yi, J., and Gu, Q. A frank-wolfe framework for efficient and effective adversarial attacks. *arXiv preprint arXiv:1811.10828*, 2018.

Cisse, M., Bojanowski, P., Grave, E., Dauphin, Y., and Usunier, N. Parseval networks: Improving robustness to adversarial examples. In *ICML*, 2017.

Danskin, J. M. *The theory of max-min and its application to weapons allocation problems*, volume 5. Springer Science & Business Media, 2012.

Das, N., Shanbhogue, M., Chen, S.-T., Hohman, F., Li, S., Chen, L., Kounavis, M. E., and Chau, D. H. Compression to the rescue: Defending from adversarial attacks across modalities. In *KDD*, 2018.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Du, S. S., Lee, J. D., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018.

Farnia, F., Zhang, J., and Tse, D. Generalizable adversarial training via spectral normalization. In *ICLR*, 2019.

Frank, M. and Wolfe, P. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3(1-2): 95–110, 1956.

Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. In *ICLR*, 2015.

Gu, S. and Rigazio, L. Towards deep neural network architectures robust to adversarial examples. *arXiv preprint arXiv:1412.5068*, 2014.

Guo, C., Rana, M., Cisse, M., and van der Maaten, L. Countering adversarial images using input transformations. In *ICLR*, 2018.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *CVPR*, 2016.

Hein, M. and Andriushchenko, M. Formal guarantees on the robustness of a classifier against adversarial manipulation. In *NeurIPS*, 2017.

Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.

Jiang, L., Meng, D., Zhao, Q., Shan, S., and Hauptmann, A. G. Self-paced curriculum learning. In *AAAI*, 2015.

Kolter, J. Z. and Wong, E. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *ICML*, 2018.

Kumar, M. P., Packer, B., and Koller, D. Self-paced learning for latent variable models. In *NeurIPS*, 2010.

Kurakin, A., Goodfellow, I., and Bengio, S. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.

Lee, J. and Raginsky, M. Minimax statistical learning with wasserstein distances. In *NeurIPS*, 2018.

Liao, F., Liang, M., Dong, Y., Pang, T., Hu, X., and Zhu, J. Defense against adversarial attacks using high-level representation guided denoiser. In *CVPR*, 2018.

Liu, Q., Liu, T., Liu, Z., Wang, Y., Jin, Y., and Wen, W. Security analysis and enhancement of model compressed deep learning systems under adversarial attacks. In *ASPDAC*, 2018.

Ma, X., Li, B., Wang, Y., Erfani, S. M., Wijewickrema, S., Schoenebeck, G., Song, D., Houle, M. E., and Bailey, J. Characterizing adversarial subspaces using local intrinsic dimensionality. *arXiv preprint arXiv:1801.02613*, 2018.

Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.

Nair, V. and Hinton, G. E. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.

Nøkland, A. Improving back-propagation by adding an adversarial gradient. *arXiv preprint arXiv:1510.04189*, 2015.

Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z. B., and Swami, A. The limitations of deep learning in adversarial settings. In *IEEE European Symposium on Security and Privacy*, 2016a.

Papernot, N., McDaniel, P., Wu, X., Jha, S., and Swami, A. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE Symposium on Security and Privacy*, 2016b.

Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z. B., and Swami, A. Practical black-box attacks against machine learning. In *Asia CCS*, 2017.

Parzen, E. On estimation of a probability density function and mode. *The Annals of Mathematical Statistics*, 33(3): 1065–1076, 1962.

Rakin, A. S., Yi, J., Gong, B., and Fan, D. Defend deep neural networks against adversarial examples via fixed and-dynamic quantized activation functions. *arXiv preprint arXiv:1807.06714*, 2018.

Ross, A. S. and Doshi-Velez, F. Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. In *AAAI*, 2018.

Samangouei, P., Kabkab, M., and Chellappa, R. Defense-GAN: Protecting classifiers against adversarial attacks using generative models. In *ICLR*, 2018.

Sinha, A., Namkoong, H., and Duchi, J. Certifiable distributional robustness with principled adversarial training. In *ICLR*, 2018.

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D., and McDaniel, P. Ensemble adversarial training: Attacks and defenses. In *ICLR*, 2018.

Xu, W., Evans, D., and Qi, Y. Feature squeezing: Detecting adversarial examples in deep neural networks. *arXiv preprint arXiv:1704.01155*, 2017.

Yan, Z., Guo, Y., and Zhang, C. Deep defense: Training dnns with improved adversarial robustness. In *NeurIPS*, 2018.

Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

Zou, D., Cao, Y., Zhou, D., and Gu, Q. Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.

# A. Proof of Theorem 1

The proof of Theorem 1 is inspired by Sinha et al. (2018). Before we prove this theorem, we need the following two technical lemmas.

**Lemma 1.** *Under Assumptions 1 and 2, we have $L_S(\boldsymbol{\theta})$ is L-smooth where $L = L_{\theta x} L_{x\theta}/\mu + L_{\theta\theta}$, i.e., for any $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$, it holds*

$$L_S(\boldsymbol{\theta}_1) \leq L_S(\boldsymbol{\theta}_2) + \langle \nabla L_S(\boldsymbol{\theta}_2), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \frac{L}{2}\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2,$$

$$\|\nabla L_S(\boldsymbol{\theta}_1) - \nabla L_S(\boldsymbol{\theta}_2)\|_2 \leq L\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2$$

*Proof.* By Assumption 2, we have for any $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2$, and $\mathbf{x}_i^*(\boldsymbol{\theta}_1), \mathbf{x}_i^*(\boldsymbol{\theta}_2)$, we have

$$f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) \leq f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2)) + \langle \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2)), \mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2) \rangle - \frac{\mu}{2}\|\mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2)\|_2^2$$

$$\leq f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2)) - \frac{\mu}{2}\|\mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2)\|_2^2, \tag{6}$$

where the inequality follows from $\langle \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2)), \mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2) \rangle \leq 0$. In addition, we have

$$f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2)) \leq f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) + \langle \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_1)), \mathbf{x}_i^*(\boldsymbol{\theta}_2) - \mathbf{x}_i^*(\boldsymbol{\theta}_1) \rangle - \frac{\mu}{2}\|\mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2)\|_2^2 \tag{7}$$

Combining (6) and (7), we obtain

$$\mu\|\mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2)\|_2^2 \leq \langle \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_1)), \mathbf{x}_i^*(\boldsymbol{\theta}_2) - \mathbf{x}_i^*(\boldsymbol{\theta}_1) \rangle$$

$$\leq \langle \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) - \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_1)), \mathbf{x}_i^*(\boldsymbol{\theta}_2) - \mathbf{x}_i^*(\boldsymbol{\theta}_1) \rangle$$

$$\leq \|\nabla_{\mathbf{x}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) - \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_1))\|_2 \|\mathbf{x}_i^*(\boldsymbol{\theta}_2) - \mathbf{x}_i^*(\boldsymbol{\theta}_1)\|_2$$

$$\leq L_{x\theta}\|\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1\|_2 \|\mathbf{x}_i^*(\boldsymbol{\theta}_2) - \mathbf{x}_i^*(\boldsymbol{\theta}_1)\|_2 \tag{8}$$

where the second inequality holds because $\langle \nabla_{\mathbf{x}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_1)), \mathbf{x}_i^*(\boldsymbol{\theta}_2) - \mathbf{x}_i^*(\boldsymbol{\theta}_1) \rangle \leq 0$, the third inequality follows from Cauchy–Schwarz inequality, and the last inequality holds due to Assumption 1. (8) immediately yields

$$\|\mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2)\|_2 \leq \frac{L_{x\theta}}{\mu}\|\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1\|_2. \tag{9}$$

Then we have for $i \in [n]$,

$$\|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2))\|_2 \leq \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_2))\|_2$$

$$+ \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_2)) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2))\|_2$$

$$\leq L_{\theta x}\|\mathbf{x}_i^*(\boldsymbol{\theta}_1) - \mathbf{x}_i^*(\boldsymbol{\theta}_2)\|_2 + L_{\theta\theta}\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2$$

$$= \left(\frac{L_{\theta x} L_{x\theta}}{\mu} + L_{\theta\theta}\right)\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2 \tag{10}$$

where the first inequality follows from triangle inequality, the second inequality holds due to Assumption 1, and the last inequality is due to (10). Finally, by the definition of $L_S(\boldsymbol{\theta})$, we have

$$\|\nabla L_S(\boldsymbol{\theta}_1) - \nabla L_S(\boldsymbol{\theta}_2)\|_2 \leq \left\|\frac{1}{n}\sum_{i=1}^n \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) - \frac{1}{n}\sum_{i=1}^n \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2))\right\|_2$$

$$\leq \frac{1}{n}\sum_{i=1}^n \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_1, \mathbf{x}_i^*(\boldsymbol{\theta}_1)) - \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_2, \mathbf{x}_i^*(\boldsymbol{\theta}_2))\|_2$$

$$\leq \left(\frac{L_{\theta x} L_{x\theta}}{\mu} + L_{\theta\theta}\right)\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2,$$

where the last inequality follows from (10). This completes the proof. $\square$

**Lemma 2.** *Under Assumptions 1 and 2, the approximate stochastic gradient $\hat{\mathbf{g}}(\theta)$ satisfies*

$$\|\hat{\mathbf{g}}(\boldsymbol{\theta}) - \mathbf{g}(\boldsymbol{\theta})\|_2 \leq L_{\theta x}\sqrt{\frac{\delta}{\mu}}. \tag{11}$$

*Proof.* We have

$$\begin{aligned}
\|\hat{\mathbf{g}}(\boldsymbol{\theta}) - \mathbf{g}(\boldsymbol{\theta})\|_2 &= \left\|\frac{1}{|\mathcal{B}|}\sum_{i\in\mathcal{B}}(\nabla_{\boldsymbol{\theta}}f(\boldsymbol{\theta},\hat{\mathbf{x}}_i(\boldsymbol{\theta})) - \nabla\bar{f}_i(\boldsymbol{\theta}))\right\|_2 \\
&\leq \frac{1}{|\mathcal{B}|}\sum_{i\in\mathcal{B}}\left\|\nabla_{\boldsymbol{\theta}}f(\boldsymbol{\theta},\hat{\mathbf{x}}_i(\boldsymbol{\theta})) - \nabla_{\boldsymbol{\theta}}f(\boldsymbol{\theta},\mathbf{x}_i^*(\boldsymbol{\theta}))\right\|_2 \\
&\leq \frac{1}{|\mathcal{B}|}\sum_{i\in\mathcal{B}}L_{\theta x}\|\hat{\mathbf{x}}_i(\boldsymbol{\theta}) - \mathbf{x}_i^*(\boldsymbol{\theta})\|_2,
\end{aligned} \tag{12}$$

where the first inequality follows from triangle inequality, and the second inequality holds due to Assumption 1. By Assumption 2, we have for any $\boldsymbol{\theta}$, and $\mathbf{x}_i^*(\boldsymbol{\theta}), \hat{\mathbf{x}}_i(\boldsymbol{\theta})$, we have

$$\mu\|\mathbf{x}_i^*(\boldsymbol{\theta}) - \hat{\mathbf{x}}_i(\boldsymbol{\theta})\|_2^2 \leq \langle\nabla_{\mathbf{x}}f(\boldsymbol{\theta},\mathbf{x}_i^*(\boldsymbol{\theta})) - \nabla_{\mathbf{x}}f(\boldsymbol{\theta},\hat{\mathbf{x}}_i(\boldsymbol{\theta})),\hat{\mathbf{x}}_i(\boldsymbol{\theta}) - \mathbf{x}_i^*(\boldsymbol{\theta})\rangle. \tag{13}$$

Since $\hat{\mathbf{x}}_i(\boldsymbol{\theta})$ is a $\delta$-approximate maximizer of $f(\boldsymbol{\theta},\hat{\mathbf{x}}_i(\boldsymbol{\theta}))$, we have

$$\langle\mathbf{x}_i^*(\boldsymbol{\theta}) - \hat{\mathbf{x}}_i(\boldsymbol{\theta}),\nabla_{\boldsymbol{\theta}}f(\boldsymbol{\theta},\hat{\mathbf{x}}_i(\boldsymbol{\theta}))\rangle \leq \delta. \tag{14}$$

In addition, we have

$$\langle\hat{\mathbf{x}}_i(\boldsymbol{\theta}) - \mathbf{x}_i^*(\boldsymbol{\theta}),\nabla_{\mathbf{x}}f(\boldsymbol{\theta},\mathbf{x}_i^*(\boldsymbol{\theta}))\rangle \leq 0. \tag{15}$$

Combining (14) and (15) gives rise to

$$\langle\hat{\mathbf{x}}_i(\boldsymbol{\theta}) - \mathbf{x}_i^*(\boldsymbol{\theta}),\nabla_{\mathbf{x}}f(\boldsymbol{\theta},\mathbf{x}_i^*(\boldsymbol{\theta})) - \nabla_{\boldsymbol{\theta}}f(\boldsymbol{\theta},\hat{\mathbf{x}}_i(\boldsymbol{\theta}))\rangle \leq \delta. \tag{16}$$

Substitute (16) into (13), we obtain

$$\mu\|\mathbf{x}_i^*(\boldsymbol{\theta}) - \hat{\mathbf{x}}_i(\boldsymbol{\theta})\|_2^2 \leq \delta,$$

which immediately yields

$$\|\mathbf{x}_i^*(\boldsymbol{\theta}) - \hat{\mathbf{x}}_i(\boldsymbol{\theta})\|_2 \leq \sqrt{\frac{\delta}{\mu}}. \tag{17}$$

Substitute (17) into (12), we obtain

$$\|\hat{\mathbf{g}}(\boldsymbol{\theta}) - \mathbf{g}(\boldsymbol{\theta})\|_2 \leq L_{\theta x}\sqrt{\frac{\delta}{\mu}},$$

which completes the proof. $\square$

Now we are ready to prove Theorem 1.

*Proof of Theorem 1.* Let $\bar{f}(\boldsymbol{\theta}) = 1/n \sum_{i=1}^{n} \min_{\mathbf{x}_i} f(\boldsymbol{\theta}, \mathbf{x}_i) = 1/n \sum_{i=1}^{n} f(\boldsymbol{\theta}, \mathbf{x}_i^*)$. By Lemma 1, we have

$$
\begin{aligned}
L_S(\boldsymbol{\theta}^{t+1}) &\leq L_S(\boldsymbol{\theta}^t) + \langle \nabla L_S(\boldsymbol{\theta}^t), \boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t \rangle + \frac{L}{2} \|\boldsymbol{\theta}^{t+1} - \boldsymbol{\theta}^t\|_2^2 \\
&= L_S(\boldsymbol{\theta}^t) - \eta_t \|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2 + \frac{L\eta_t^2}{2} \|\hat{\mathbf{g}}(\boldsymbol{\theta}^t)\|_2^2 + \eta_t \langle \nabla L_S(\boldsymbol{\theta}^t), \nabla L_S(\boldsymbol{\theta}^t) - \hat{\mathbf{g}}(\boldsymbol{\theta}^t) \rangle \\
&= L_S(\boldsymbol{\theta}^t) - \eta_t \left(1 - \frac{L\eta_t}{2}\right) \|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2 + \eta_t \left(1 - L\eta_t\right) \langle \nabla L_S(\boldsymbol{\theta}^t), \nabla L_S(\boldsymbol{\theta}^t) - \hat{\mathbf{g}}(\boldsymbol{\theta}^t) \rangle \\
&\quad + \frac{L\eta_t^2}{2} \|\hat{\mathbf{g}}(\boldsymbol{\theta}^t) - \nabla L_S(\boldsymbol{\theta}^t)\|_2^2 \\
&= L_S(\boldsymbol{\theta}^t) - \eta_t \left(1 - \frac{L\eta_t}{2}\right) \|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2 + \eta_t \left(1 - L\eta_t\right) \langle \nabla L_S(\boldsymbol{\theta}^t), \mathbf{g}(\boldsymbol{\theta}^t) - \hat{\mathbf{g}}(\boldsymbol{\theta}^t) \rangle \\
&\quad + \eta_t \left(1 - L\eta_t\right) \langle \nabla L_S(\boldsymbol{\theta}^t), \nabla L_S(\boldsymbol{\theta}^t) - \mathbf{g}(\boldsymbol{\theta}^t) \rangle + \frac{L\eta_t^2}{2} \|\hat{\mathbf{g}}(\boldsymbol{\theta}^t) - \mathbf{g}(\boldsymbol{\theta}^t) + \mathbf{g}(\boldsymbol{\theta}^t) - \nabla L_S(\boldsymbol{\theta}^t)\|_2^2 \\
&\leq L_S(\boldsymbol{\theta}^t) - \frac{\eta_t}{2} \|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2 + \frac{\eta_t}{2} \left(1 - L\eta_t\right) \|\hat{\mathbf{g}}(\boldsymbol{\theta}^t) - \mathbf{g}(\boldsymbol{\theta}^t)\|_2^2 \\
&\quad + \eta_t \left(1 - L\eta_t\right) \langle \nabla L_S(\boldsymbol{\theta}^t), \nabla L_S(\boldsymbol{\theta}^t) - \mathbf{g}(\boldsymbol{\theta}^t) \rangle + L\eta_t^2 \left(\|\hat{\mathbf{g}}(\boldsymbol{\theta}^t) - \mathbf{g}(\boldsymbol{\theta}^t)\|_2^2 + \|\mathbf{g}(\boldsymbol{\theta}^t) - \nabla L_S(\boldsymbol{\theta}^t)\|_2^2\right),
\end{aligned}
$$

where the last inequality is due to the Young's inequality. Note that we have $\eta_t \leq 1/L$ because we choose $\eta_t = \eta = \sqrt{\Delta/(TL\sigma^2)}$ and $T \geq (\Delta L)/\sigma^2$. Taking expectation on both sides of the above inequality conditioned on $\boldsymbol{\theta}^t$, we have

$$
\begin{aligned}
\mathbb{E}[L_S(\boldsymbol{\theta}^{t+1}) - L_S(\boldsymbol{\theta}^t)|\boldsymbol{\theta}^t] &\leq -\frac{\eta_t}{2} \|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2 + \frac{\eta_t}{2} \left(1 + L\eta_t\right) \frac{L_{\theta x}^2 \delta}{\mu} + L\eta_t^2 \sigma^2 \\
&\leq -\frac{\eta_t}{2} \|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2 + \eta_t \frac{L_{\theta x}^2 \delta}{\mu} + L\eta_t^2 \sigma^2,
\end{aligned} \tag{18}
$$

where the first inequality uses the fact that $\mathbb{E}[\mathbf{g}(\boldsymbol{\theta}^t)] = \nabla L_S(\boldsymbol{\theta}^t)$, Assumption 3, and Lemma 2, and the second inequality uses the fact that $\eta_t \leq 1/L$. Taking telescope sum of (18) over $t = 0, \dots, T-1$, we obtain that

$$
\sum_{t=0}^{T-1} \frac{\eta_t}{2} \mathbb{E}\left[\|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2\right] \leq \mathbb{E}[L_S(\boldsymbol{\theta}^0) - L_S(\boldsymbol{\theta}^T)] + \sum_{t=0}^{T-1} \eta_t \frac{L_{\theta x}^2 \delta}{\mu} + L \sum_{t=0}^{T-1} \eta_t^2 \sigma^2.
$$

Recall that $\eta_t = \eta = \sqrt{\Delta/(TL\sigma^2)}$ where $L = L_{\theta x} L_{x\theta}/\mu + L_{\theta\theta}$, we can show that

$$
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[\|\nabla L_S(\boldsymbol{\theta}^t)\|_2^2\right] \leq 4\sigma \sqrt{\frac{L\Delta}{T}} + \frac{2L_{\theta x}^2 \delta}{\mu}.
$$

This completes the proof. $\qquad\square$