# UNDERSTANDING AND QUANTIFYING ADVERSARIAL EXAMPLES EXISTENCE IN LINEAR CLASSIFICATION

**Xupeng Shi**
Department of Mathematics
Northeastern University
Boston, USA
shi.xup@husky.neu.edu

**A. Adam Ding**
Department of Mathematics
Northeastern University
Boston, USA
a.ding@northeastern.edu

October 29, 2019

## ABSTRACT

State-of-art deep neural networks (DNN) are vulnerable to attacks by adversarial examples: a carefully designed small perturbation to the input, that is imperceptible to human, can mislead DNN. To understand the root cause of adversarial examples, we quantify the probability of adversarial example existence for linear classifiers. Previous mathematical definition of adversarial examples only involves the overall perturbation amount, and we propose a more practical relevant definition of strong adversarial examples that separately limits the perturbation along the signal direction also. We show that linear classifiers can be made robust to strong adversarial examples attack in cases where no adversarial robust linear classifiers exist under the previous definition. The quantitative formulas are confirmed by numerical experiments using a linear support vector machine (SVM) classifier. The results suggest that designing general strong-adversarial-robust learning systems is feasible but only through incorporating human knowledge of the underlying classification problem.

## 1 Introduction

The deep neural networks (DNN) are widely used as the state-of-art machining learning classification systems due to its great performance gains in recent years. Meanwhile adversarial examples, first pointed out by Szegedy et al. (2014), emerges as a novel peculiar security threat against such systems: a small perturbation that is unnoticeable to human eyes can cause the DNNs to misclassify. Various adversarial algorithms have since been developed to efficiently find adversarial examples (Goodfellow et al., 2015; Moosavi-Dezfooli et al., 2016; Carlini and Wagner, 2017; Madry et al., 2018). The adversarial examples have also been demonstrated to misled DNN based classification systems in physical world applications (Sharif et al., 2016; Brown et al., 2017; Kurakin et al., 2018; Athalye et al., 2018a). Various defense methods have also been proposed to prevent adversarial example attacks: Adversarial training (Szegedy et al., 2014; Goodfellow et al., 2015); Defensive distillation Papernot et al. (2016); Minmax robust training (Madry et al., 2018; Sinha et al., 2018); Input transformation Xu et al. (2017). However, many of the defenses are shown to be vulnerable to attacks taking such defense strategies into consideration (Athalye et al., 2018b).

Recently, Shafahi et al. (2019) showed that, for two classes of data distributed with bounded probability densities on a compact region of a high dimensional space, no classifier can both have low misclassification rate and be robust to adversarial examples attack. So are we left hopeless against such threat? Theoretical analysis for understanding adversarial examples is needed to address this issue. Goodfellow et al. (2015); Fawzi et al. (2018) pointed out that susceptibility of DNN classifiers to adversarial attacks could be related to their locally linear behaviours. The existence of adversarial examples is not unique to DNN, traditional linear classifiers also have adversarial examples. In this paper, we extend the understanding of adversarial examples by quantifying the probability of their existence for a simple case of linear classifiers that performs binary classification on Gaussian mixture data.

In previous literature, a data point $x$ is mathematically defined as having an adversarial example $x' = x + v$ when the perturbation amount $\|v\|$ is small and $x'$ is classified differently from $x$. This definition does not exclude genuine

signal perturbation. For example, if a dog image $\boldsymbol{x}$ is perturbed to an image $\boldsymbol{x}'$ that is classified as a cat by both human and the machine classifier, then $\boldsymbol{x}'$ should not be an adversarial example even if $\|\boldsymbol{v}\| = \|\boldsymbol{x}' - \boldsymbol{x}\|$ is small. The proper definition needs to capture the novelty of adversarial examples attack: while a human would consider two images $\boldsymbol{x}'$ and $\boldsymbol{x}$ very similar and consider both clearly as dogs, a machine classifier misclassifies $\boldsymbol{x}'$ as a cat. While defining genuine signal perturbation for general learning problems is difficult mathematically, the signal perturbation is clear in the binary linear classification for Gaussian mixture data. We therefore propose a new definition of strong-adversarial examples that limits the perturbation amount in the signal direction separately from the limit on overall perturbation amount.

In this paper, we derive quantitative formulas for the probabilities of adversarial and strong-adversarial examples existence in the binary linear classification problem. Our quantitative analysis shows that an adversarial-robust linear classifier requires much higher signal-to-noise ratio (SNR) in data than a good performing classifier does. Therefore, in many practical applications, adversarial-robust classifiers may not be available nor are such classifiers desirable. On the contrary, useful strong-adversarial-robust linear classifiers exists at the SNR similar to that required by the existence of any useful linear classifiers, however, they require better designed training algorithms.

The paper is organized as follows. Section 2 presents the notations and definitions of (strong-)adversarial examples and derive explicit formulas for the probability of their existence. Section 3 presents numerical experiments. The formulas are confirmed experimentally, and then are used to illustrate their implication on the vulnerability against (strong-)adversarial example attacks. Section 4 discusses how our results relate to some works in literature and summarize their implication on general adversarial attack defenses.

## 2   Adversarial Rates Analysis of Linear Binary Classifier on Gaussian Mixture Data

We first introduce our definitions of adversarial and strong-adversarial examples, and then we characterize their existence through defining sets. Using the defining sets, we derive explicit probability rates of (strong-)adversarial examples existence for linear classifiers on Gaussian mixture data.

### 2.1   Definition of Adversarial and Strong-Adversarial Examples

The classical adversarial examples are defined as follows:

**Definition 1.** [1] *Given a classifier $C$, an $\varepsilon$-adversarial example of a data vector $\boldsymbol{x}$ is another data vector $\boldsymbol{x}'$ such that $\|\boldsymbol{x} - \boldsymbol{x}'\| \le \varepsilon$ but $C(\boldsymbol{x}) \neq C(\boldsymbol{x}')$.*

Without loss of generality, in this paper we focus on $\ell_2$ norm perturbations. If not specified, $\|\cdot\|$ in the following refers to the $\ell_2$ norm. The general $\ell_p$ norm ($p \ge 1$) perturbation is studied in the Appendix 5, and the results will be stated in the discussion section.

For a general machine classification problem, it is reasonable to only consider adversarial examples since the signal direction is often not easily definable mathematically. Here we consider the simple binary linear classification of Gaussian mixture data where the signal direction can be clearly distinguished. For two classes labeled '+' and '−' respectively, a linear classifier is $C(\boldsymbol{x}; \boldsymbol{w}, b) = \{\boldsymbol{w} \cdot \boldsymbol{x} + b > 0\}$ where '·' denotes the inner product of two vectors. Here the parameters $\boldsymbol{w}$ and $b$ are respectively the weight vector and the bias term. For the classical Gaussian mixture data problem, for each of the two classes, the $d$-dimensional data vector $\boldsymbol{x}$ comes from a multivariate Gaussian distribution $N(\boldsymbol{\mu}_i, \sigma_i^2 \boldsymbol{I}_d)$, $i = $ '+' or '−'. Notice the optimal ideal classifier here is the Bayes classifier $C(\boldsymbol{x}; \boldsymbol{\mu}, \bar{\boldsymbol{\mu}}) = \{\boldsymbol{\mu} \cdot (\boldsymbol{x} - \bar{\boldsymbol{\mu}}) > 0\}$[2] where $\boldsymbol{\mu} = \frac{1}{2}(\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-), \bar{\boldsymbol{\mu}} = \frac{1}{2}(\boldsymbol{\mu}_+ + \boldsymbol{\mu}_-)$.

For this problem, the data distributions of the two classes only differ in their means $\boldsymbol{\mu}_+$ and $\boldsymbol{\mu}_-$. Thus the signal direction is $\boldsymbol{\mu}_0 = \boldsymbol{\mu} / \|\boldsymbol{\mu}\|$. Adding $2 \|\boldsymbol{\mu}\|$ amount of perturbation along the signal direction changes the '−' class data distribution to the '+' class data distribution exactly, rending all classifiers unable to defend against such a perturbation.

In previous literature, the adversarial examples definition does not limit perturbation along the signal direction, therefore we propose a new definition that limits the perturbation along the signal direction separately by an amount $\delta$, we will refer these examples as *strong-adversarial examples* .

**Definition 2.** *Given a classifier $C$, an $(\varepsilon, \delta)$-strong-adversarial example of a data vector $\boldsymbol{x}$ is another data vector $\boldsymbol{x}'$ such that $\|\boldsymbol{x} - \boldsymbol{x}'\| \le \varepsilon$ and $|(\boldsymbol{x} - \boldsymbol{x}') \cdot \boldsymbol{\mu}_0| \le \delta$ but $C(\boldsymbol{x}) \neq C(\boldsymbol{x}')$.*

---

[1] We don't distinguish the targeted and untargeted adversarial examples here because for binary classification they are the same.

[2] Here we just use the optimal Bayes classfier for balanced case since we are focusing on the balanced case in the following text.

To illustrate the difference between the adversarial examples and the strong-adversarial examples, we consider the following examples visualized in Figure 1. Here, Figure 1(a) shows a data vector $\boldsymbol{x}$ of dimension $d = 19 \times 19 = 361$ from the '+' class. To visualize, each component of the data vector is mapped onto $[0, 1]$ via function $\frac{1}{2}(\tanh \frac{2x}{3} + 1)$ and then displayed in grey scale as a $19 \times 19$ image (Carlini and Wagner, 2017).
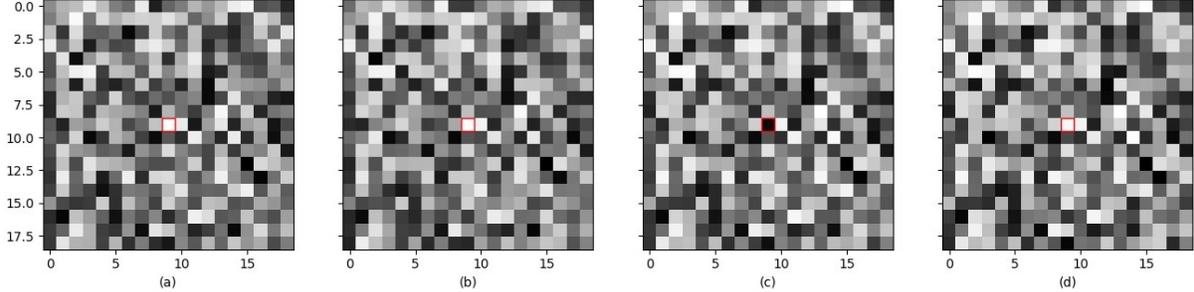


Figure 1: (a) a data point $\boldsymbol{x}$ from the '+' class; (b) a randomly perturbed $\boldsymbol{x}'$; (c) an adversarial $\boldsymbol{x}'$ but not strong-adversarial; (d) a strong-adversarial $\boldsymbol{x}'$. All three perturbations are of the same amount $\varepsilon = 5.7$ and $\|\boldsymbol{\mu}\| = 4$. The center grid cell within the red boundary contains the real class signal.

The two means $\boldsymbol{\mu}_+$ and $\boldsymbol{\mu}_-$ are chosen to be zero at every component of the vector except the component corresponding to center grid cell (shown with red boundary in Figure 1). Hence the optimal Bayes classifier identifies the image as from '+' (or '−') class when the center grid cell within the red boundary appears to be white (or black). With a perturbation amount of $\varepsilon = 0.3 \times 19 = 5.7$, Figure 1(b) shows a randomly perturbed $\boldsymbol{x}'$ which is hardly distinguishable from the first image $\boldsymbol{x}$ to the human eye. This confirms that, in defending against realistic threats, $\varepsilon$ of magnitude $O(\sqrt{d})$ needs to be studied. (Detailed discussion of $\varepsilon$ order is in subsection 2.3.)

For a trained support vector machine (SVM) classifier, Figure 1(c) and (d) shows two adversarial examples with the same $\varepsilon = 5.7$, but only the last one in (d) is strong-adversarial for $\delta = 1.2$. (Section 3 provides detailed setup of this experiment.) The adversarial attacks present a novel threat: a machine classifier misclassifies the perturbed data points that a human would not have noted the difference. We can see that our strong-adversarial example definition focus attention on this novel threat. In contrast, under the traditional definition, the adversarial examples include examples similar to Figure 1(c) that would indeed be classified by human into another class. We now quantitatively analyze the existence of adversarial and strong-adversarial examples.

## 2.2   The Defining Sets

Here we characterize the defining sets where the (strong-)adversarial examples exist. Then we quantify the probability of data falling into these defining sets in the next subsection 2.3.

We denote $\Omega_\varepsilon = \{\boldsymbol{x} : \boldsymbol{x} \text{ has an } \varepsilon\text{-adversarial example}\}$ and $\Omega_{\varepsilon,\delta} = \{\boldsymbol{x} : \boldsymbol{x} \text{ has an } (\varepsilon, \delta)\text{-strong-adversarial example}\}$. Furthermore, for a fixed perturbation $\boldsymbol{v}$, we denote the set where $\boldsymbol{v}$ changes classification as $\Omega(\boldsymbol{v}) = \{\boldsymbol{x} \in \mathbb{R}^d : C(\boldsymbol{x} + \boldsymbol{v}) \neq C(\boldsymbol{x})\}$.

For any data point $\boldsymbol{x}$ in $\Omega_\varepsilon$, there exists a $\boldsymbol{v}$ with $\|\boldsymbol{v}\| \leq \varepsilon$ such that $\boldsymbol{x} + \boldsymbol{v}$ is classified differently from $\boldsymbol{x}$. In other words, the distance of $\boldsymbol{x}$ from the classifier's decision boundary is less than $\varepsilon$. For a linear classifier $C(\boldsymbol{x}; \boldsymbol{w}, b) = \{\boldsymbol{w} \cdot \boldsymbol{x} + b > 0\}$, the normal direction of its decision boundary is $\boldsymbol{v}_0 = \boldsymbol{w}/\|\boldsymbol{w}\|$. Thus, perturbing $\boldsymbol{x}$ by $\varepsilon$ amount along one of the two directions $\boldsymbol{v}_0$ or $-\boldsymbol{v}_0$ will cross the linear decision boundary. That is, $\Omega_\varepsilon \subseteq \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0)$. Since it is obvious from the definition that $\Omega_\varepsilon = \bigcup_{\|\boldsymbol{v}\| \leq \varepsilon} \Omega(\boldsymbol{v}) \supseteq \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0)$, we have $\Omega_\varepsilon = \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0)$. In summary, to judge if $\boldsymbol{x} \in \Omega_\varepsilon$, we only need to check the perturbation along the normal direction $\boldsymbol{v}_0$.

In contrast, our definition of strong-adversarial examples only allows $\delta$ amount of perturbation along the signal notation $\boldsymbol{\mu}_0$, hence it is not sufficient to only check perturbations $\varepsilon\boldsymbol{v}_0$ and $-\varepsilon\boldsymbol{v}_0$ for judging if $\boldsymbol{x} \in \Omega_{\varepsilon,\delta}$. Let $\theta$ denote the deflected angle between $\boldsymbol{\mu}_0$ and $\boldsymbol{v}_0$. (Without loss of generality, we choose the $\theta$ value such that $0 \leq \theta \leq \pi/2$.) Then we can decompose $\boldsymbol{v}_0$ into two components along and orthogonal to the signal direction $\boldsymbol{\mu}_0$ respectively. That is, $\boldsymbol{v}_0 = \cos\theta\boldsymbol{\mu}_0 + \sin\theta\boldsymbol{n}_0$ where $\boldsymbol{n} = \boldsymbol{v}_0 - (\boldsymbol{v}_0 \cdot \boldsymbol{\mu}_0)\boldsymbol{\mu}_0$ and $\boldsymbol{n}_0 = \boldsymbol{n}/\|\boldsymbol{n}\|$. When $\varepsilon\cos\theta \leq \delta$, the adversarial example resulting from the $\varepsilon\boldsymbol{v}_0$ perturbation is also strong-adversarial by definition. When $\varepsilon\cos\theta > \delta$, however, $\varepsilon\boldsymbol{v}_0$ is no longer an allowable perturbation in the strong-adversarial example definition. Then we need to check whether classification

change is caused by a perturbation of $\delta$ amount along $\boldsymbol{\mu}_0$ direction and $\sqrt{\varepsilon^2 - \delta^2}$ amount along $\boldsymbol{n}_0$ direction. That is, , to judge if $\boldsymbol{x} \in \Omega_{\varepsilon,\delta}$, we need to check perturbations $\boldsymbol{u}_2 = \delta\boldsymbol{\mu}_0 + \sqrt{\varepsilon^2 - \delta^2}\boldsymbol{n}_0$ and $-\boldsymbol{u}_2$. We summarize the defining sets characterization in the following lemma whose detailed proof is in the Appendix 5.1.

**Lemma 1.** *The defining sets for $\varepsilon$-adversarial and $(\varepsilon, \delta)$-strong-adversarial examples are given by:*

$$\Omega_\varepsilon = \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0); \quad \Omega_{\varepsilon,\delta} = \Omega(\boldsymbol{u}_2) \cup \Omega(-\boldsymbol{u}_2) \tag{1}$$

*where $\boldsymbol{u}_2 = \beta\boldsymbol{\mu}_0 + \sqrt{\varepsilon^2 - \beta^2}\boldsymbol{n}_0, \beta = \min(\varepsilon\cos\theta, \delta)$.*

Next, we use these defining sets to quantify the probabilities of (strong-)adversarial example existence.

## 2.3    Adversarial and Strong-Adversarial Rates

For the binary classification problem, a random data vector comes from the Gaussian mixture distribution $p(\boldsymbol{x}) = \lambda_+\varphi_+(\boldsymbol{x}) + \lambda_-\varphi_-(\boldsymbol{x})$, where $\varphi_i(\boldsymbol{x})$ is the probability density function of the multivariate Gaussian $N(\boldsymbol{\mu}_i, \sigma_i^2\boldsymbol{I}_d)$ and $\lambda_i$ is the probability that the data vector belongs to the class of $i = $ '$+$' or '$-$'. For simplicity, we focus on the balanced classes case of $\lambda_+ = \lambda_- = 0.5$ and also $\sigma_+ = \sigma_- = \sigma$.

**Adversarial Rate**    For a random data vector $\boldsymbol{x}$ from the '$+$' class, it has an $\varepsilon$-adversarial example $\boldsymbol{x}'$ if it is classified correctly by $\boldsymbol{w} \cdot \boldsymbol{x} + b > 0$ and $\boldsymbol{x} \in \Omega(-\varepsilon\boldsymbol{v}_0)$. Thus the adversarial rate from the '$+$' class is

$$\lambda_+ pr[\boldsymbol{w} \cdot \boldsymbol{x} + b > 0, \boldsymbol{w} \cdot (\boldsymbol{x} - \varepsilon\boldsymbol{v}_0) + b < 0 \,|\varphi_+(\boldsymbol{x})] = 0.5 pr[0 < \boldsymbol{w} \cdot \boldsymbol{x} + b < \varepsilon\|\boldsymbol{w}\| \,|\varphi_+(\boldsymbol{x})]. \tag{2}$$

Since under the multivariate Gaussian $N(\boldsymbol{\mu}_+, \sigma^2\boldsymbol{I}_d)$ distribution $\varphi_+(\boldsymbol{x})$, $\boldsymbol{w} \cdot \boldsymbol{x} + b$ is a univariate Gaussian random variable with mean $\boldsymbol{w} \cdot \boldsymbol{\mu}_+ + b$ and variance $\|\boldsymbol{w}\|^2 \sigma^2$, the above quantity becomes

$$0.5\left[\Phi\left(\frac{\varepsilon\|\boldsymbol{w}\| - (\boldsymbol{w} \cdot \boldsymbol{\mu}_+ + b)}{\|\boldsymbol{w}\|\sigma}\right) - \Phi\left(\frac{-(\boldsymbol{w} \cdot \boldsymbol{\mu}_+ + b)}{\|\boldsymbol{w}\|\sigma}\right)\right]. \tag{3}$$

Here $\Phi(\cdot)$ denotes the cumulative distribution function (CDF) of the standard Gaussian distribution $N(0,1)$. Similarly, the adversarial rate from the '$-$' class is

$$\lambda_- pr[-\varepsilon\|\boldsymbol{w}\| < \boldsymbol{w} \cdot \boldsymbol{x} + b < 0|\varphi_-(\boldsymbol{x})] = 0.5\left[\Phi\left(\frac{-(\boldsymbol{w} \cdot \boldsymbol{\mu}_- + b)}{\|\boldsymbol{w}\|\sigma}\right) - \Phi\left(\frac{-\varepsilon\|\boldsymbol{w}\| - (\boldsymbol{w} \cdot \boldsymbol{\mu}_- + b)}{\|\boldsymbol{w}\|\sigma}\right)\right]. \tag{4}$$

Recall $\boldsymbol{\mu} = \frac{1}{2}(\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-), \bar{\boldsymbol{\mu}} = \frac{1}{2}(\boldsymbol{\mu}_+ + \boldsymbol{\mu}_-)$. If we denote $b' = \boldsymbol{w} \cdot \bar{\boldsymbol{\mu}} + b$, then we can rewritten the expressions as $\boldsymbol{w} \cdot \boldsymbol{\mu}_\pm + b = \pm\boldsymbol{w} \cdot \boldsymbol{\mu} + b'$. Combining equations (3) and (4), we have the overall adversarial rate as

$$\begin{aligned} p_{adv} &= 0.5\left[\Phi\left(\frac{\varepsilon}{\sigma} - \frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|\sigma}\right) - \Phi\left(-\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|\sigma}\right) + \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|\sigma}\right) - \Phi\left(-\frac{\varepsilon}{\sigma} + \frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|\sigma}\right)\right] \\ &= 0.5\left[\Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|\sigma}\right) - \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|\sigma} - \frac{\varepsilon}{\sigma}\right) + \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|\sigma}\right) - \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|\sigma} - \frac{\varepsilon}{\sigma}\right)\right] \end{aligned} \tag{5}$$

Also notice that the misclassification rates from the two classes are respectively $\lambda_+\Phi[-(\boldsymbol{w} \cdot \boldsymbol{\mu}_+ + b)/(\|\boldsymbol{w}\|\sigma)] = 0.5\{1 - \Phi[(\boldsymbol{w} \cdot \boldsymbol{\mu} + b')/(\|\boldsymbol{w}\|\sigma)]\}$ and $\lambda_-\{1 - \Phi[-(\boldsymbol{w} \cdot \boldsymbol{\mu}_- + b)/(\|\boldsymbol{w}\|\sigma)]\} = 0.5\{1 - \Phi[(\boldsymbol{w} \cdot \boldsymbol{\mu} - b')/(\|\boldsymbol{w}\|\sigma)]\}$. Thus the overall misclassification rate is

$$p_m = 1 - 0.5\left[\Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|\sigma}\right) + \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|\sigma}\right)\right]. \tag{6}$$

We combine equations (5) and (6) into the following Theorem.

**Theorem 1.** *The overall adversarial rate of a linear classifier for the balanced Gaussian mixture data is*

$$p_{adv} = 1 - p_m - 0.5\left[\Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|\sigma} - \frac{\varepsilon}{\sigma}\right) + \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|\sigma} - \frac{\varepsilon}{\sigma}\right)\right]. \tag{7}$$

To be robust against adversarial attacks, a linear classifier needs a low adversarial rate. For the classifier to be useful, it also needs a low misclassification rate. Hence we should look at the sum of misclassification rate and adversarial rate, which we call the *adversarial-error* rate:

$$p_{err} = p_{adv} + p_m = 1 - 0.5\left[\Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|\sigma} - \frac{\varepsilon}{\sigma}\right) + \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|\sigma} - \frac{\varepsilon}{\sigma}\right)\right] \tag{8}$$

Comparing equation (8) with (6), we can see why adversarial-robustness is hard to achieve.

First, the misclassification rate $p_m$ in (6) is minimized by the Bayes classifier with $b' = 0$ and $\boldsymbol{w} \cdot \boldsymbol{\mu} = \|\boldsymbol{w}\| \|\boldsymbol{\mu}\|$. Hence the best $p_m$ value is $1 - \Phi(\|\boldsymbol{\mu}\| /\sigma)$. There exists useful classifiers when $\|\boldsymbol{\mu}\| /\sigma$ is big enough to make $1 - \Phi(\|\boldsymbol{\mu}\| /\sigma)$ small. This is achieved for $\|\boldsymbol{\mu}\| /\sigma = O(1)$. For example, when $\|\boldsymbol{\mu}\| /\sigma = 3$, the misclassification rate of the Bayes classifier is around $0.1\%$.

However, to achieve a low adversarial-error rate in (8), the required SNR $\|\boldsymbol{\mu}\| /\sigma$ can be much bigger. When $\boldsymbol{w} \cdot \boldsymbol{\mu} > \varepsilon \|\boldsymbol{w}\|$, a lower bound for the adversarial-error rate is

$$p_{err} \geq 1 - \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu}}{\|\boldsymbol{w}\| \sigma} - \frac{\varepsilon}{\sigma} \right) \geq 1 - \Phi\left( \frac{\|\boldsymbol{\mu}\|}{\sigma} - \frac{\varepsilon}{\sigma} \right). \tag{9}$$

Therefore, the existence of a useful adversarial-robust linear classifier requires $\|\boldsymbol{\mu}\| /\sigma - \varepsilon/\sigma = O(1)$ instead. Notice that, for this Gaussian mixture data setup, the noise in each class follows the $N(\boldsymbol{0}, \sigma^2 \boldsymbol{I}_d)$ distribution with an expected square of $\ell_2$ norm of $d\sigma^2$. Therefore, for a positive constant value $\eta_a < 1$, the perturbation amount of $\varepsilon = \eta_a \sqrt{d}\sigma$ is smaller than the average noise in data and generally is hard to detect. Hence, for the typical high-dimensional data applications, an adversarial-robust linear classifier needs to protect against perturbation amount of $\varepsilon = O(\sqrt{d})$ which implies that $\|\boldsymbol{\mu}\| /\sigma = O(\sqrt{d})$ is needed from equation (9). Next, we show that this high SNR requirement is not needed for a strong-adversarial-robust linear classifier.

**Strong-Adversarial Rate**    The derivation of the strong-adversarial rate is very similar to that of the adversarial rate. From equation (1), the difference between the adversarial defining set and the strong-adversarial defining set is only that $\varepsilon \boldsymbol{v}_0$ is replaced by $\boldsymbol{u}_2 = \beta \boldsymbol{\mu}_0 + \sqrt{\varepsilon^2 - \beta^2} \boldsymbol{n}_0$. Hence the strong-adversarial rate from the '+' class is

$$0.5 pr[0 < \boldsymbol{w} \cdot \boldsymbol{x} + b < \boldsymbol{w} \cdot \boldsymbol{u}_2 | \varphi_+(\boldsymbol{x})].$$

Since $\boldsymbol{w} \cdot \boldsymbol{\mu}_0 = \|\boldsymbol{w}\| \cos\theta$ and $\boldsymbol{w} \cdot \boldsymbol{n}_0 = \|\boldsymbol{w}\| \sin\theta$, we have $\boldsymbol{w} \cdot \boldsymbol{u}_2 = (\beta \cos\theta + \sqrt{\varepsilon^2 - \beta^2} \sin\theta) \|\boldsymbol{w}\|$ where $\beta = \min(\varepsilon \cos\theta, \delta)$. We denote

$$g(\varepsilon, \delta, \theta) = \beta \cos\theta + \sqrt{\varepsilon^2 - \beta^2} \sin\theta. \tag{10}$$

Thus replacing $\varepsilon \|\boldsymbol{w}\|$ by $g(\varepsilon, \delta, \theta) \|\boldsymbol{w}\|$ in equations from (3) to (8), we have the following Theorem.

**Theorem 2.** *The overall strong adversarial rate and strong-adversarial-error rate of a linear classifier are*

$$p_{s-adv} = 1 - p_m - 0.5 \left[ \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} + b'}{\|\boldsymbol{w}\| \sigma} - \frac{g(\varepsilon, \delta, \theta)}{\sigma} \right) + \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} - b'}{\|\boldsymbol{w}\| \sigma} - \frac{g(\varepsilon, \delta, \theta)}{\sigma} \right) \right] \tag{11}$$

$$p_{s-err} = p_{s-adv} + p_m = 1 - 0.5 \left[ \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} + b'}{\|\boldsymbol{w}\| \sigma} - \frac{g(\varepsilon, \delta, \theta)}{\sigma} \right) + \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} - b'}{\|\boldsymbol{w}\| \sigma} - \frac{g(\varepsilon, \delta, \theta)}{\sigma} \right) \right] \tag{12}$$

Compared to the analysis above, the existence of a useful strong-adversarial-robust linear classifier requires $\|\boldsymbol{\mu}\| /\sigma - g(\varepsilon, \delta, \theta)/\sigma = O(1)$ instead. Besides the overall perturbation amount $\varepsilon$, the function $g(\varepsilon, \delta, \theta)$ in equation (10) is also affected by two other factors: the signal direction perturbation amount $\delta$ and the angle $\theta$ between the classifier and the ideal Bayes classifier. What is the practical relevant amount $\delta$ we should study? Let $\delta = \eta_s \mu = \eta_s \|\boldsymbol{\mu}\|$. When $\eta_s > 1$, a $\delta$ amount perturbation along the signal direction to all '+' class data points will make more than half of them be classified as '−' by the Bayes classifier (also to human eye, e.g., Figure 1(c)). Therefore, when studying real strong-adversarial perturbations (imperceptible to human but confuses machine) mathematically, we need to focus on $\eta_s < 1$. That is, $\delta = O(1)$. Compared to the overall perturbation amount $\varepsilon = O(\sqrt{d})$ discussed earlier, we see that $\delta \ll \varepsilon$ for typical high-dimensional data applications. When $\delta \ll \varepsilon$, $g(\varepsilon, \delta, \theta) \approx \delta \cos\theta + \varepsilon \sin\theta$. Hence if the linear classifier is well-trained to have small $\theta$ and small bias $b'$ (i.e., very close to the Bayes classifier), then its strong-adversarial-error rate is approximately $1 - \Phi[(1 - \eta_s)\|\boldsymbol{\mu}\|/\sigma]$, which can be made small when SNR $\|\boldsymbol{\mu}\| /\sigma$ is of order $O(1)$. That is, with good training, we can find a useful strong-adversarial-robust linear classifier when $\|\boldsymbol{\mu}\| /\sigma = O(1)$. In contrast, no training can make the linear classifier to be useful and adversarial-robust unless the SNR $\|\boldsymbol{\mu}\| /\sigma$ is much bigger, at the order of $O(\sqrt{d})$.

The conclusion for the analysis using $\ell_p$ norm (see Appendix 5 for details) is similar. There exists a useful strong-adversarial-robust linear classifier for constant order SNR $\|\boldsymbol{\mu}\| /\sigma = O(1)$, but a useful $\ell_p$-adversarial-robust linear classifier only exists when SNR is much bigger, at the order of $O(d^{min(1/p, 1/2)})$.

## 3   Numerical Studies and Analysis of Adversarial Examples

### 3.1   (Strong-)Adversarial Rates for the Linear SVM

**Settings**   We first conduct numerical experiments of a support vector machine (SVM) classifier on the Gaussian mixture data. We randomly generate 5000 data points from the balanced mixture distribution $0.5N(\boldsymbol{\mu}_+, \sigma^2 \boldsymbol{I}_d) + 0.5N(\boldsymbol{\mu}_-, \sigma^2 \boldsymbol{I}_d)$, and randomly splits them into 4000 train data and 1000 test data. We set $\boldsymbol{\mu}_+ = -\boldsymbol{\mu}_- = [\mu, 0, \cdots, 0]$, $\sigma = 1$ and $d = 19 \times 19$. A linear SVM is trained on the training data using the python *scikit-learn* package and its default setting. Then for each test data vector, we check if it has any adversarial and strong-adversarial example, for $\varepsilon = \eta_a \sqrt{d}\sigma = 19\eta_a$ and $\delta = \eta_s \mu$. Figure 1 earlier visualizes one such test data vector and its adversarial and strong-adversarial examples for $\eta_a = \eta_s = 0.3$. We conduct this experiment for various values of $\eta = \eta_a = \eta_s$ and $\mu$, and for each parameter combination, the simulation is repeated 1000 times. Figure 2 plots three empirical rates (misclassification, adversarial-error and strong-adversarial-error), each averaged over the 1000 simulations, against $\mu$ values, together with corresponding quantitative formulas from equations (6), (8) and (12). Figure 2(a)-(c) shows the results for three perturbation levels of $\eta = 0.05, 0.1, 0.3$, with the empirical quantities shown with symbols and the quantitative formulas shown in curves. The plots show very good agreement between the formulas with actual empirical proportions.
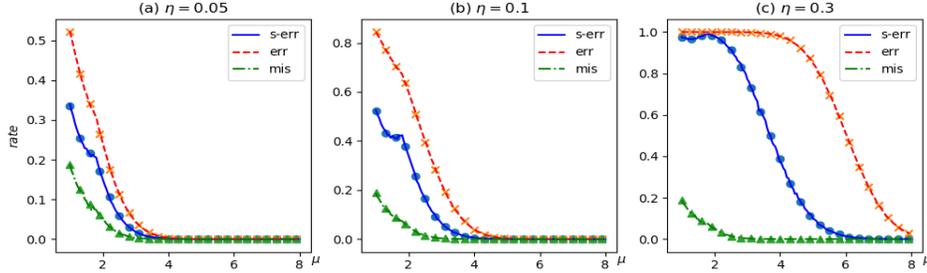


Figure 2: Empirical probabilities and their theoretical values calculated from equations (6), (8) and (12), plotted versus $\mu$. (a) $\eta = 0.05$, (b) $\eta = 0.1$, (c) $\eta = 0.3$

In our simulation, $\mu = \|\boldsymbol{\mu}\| / 1 = \|\boldsymbol{\mu}\| / \sigma$ is the SNR. Figure 2 shows that SVM have pretty good performance in terms of misclassification rate once the SNR exceeds 2. However, it is not robust to (strong-)adversarial attacks when $\mu = 2$, and will only become robust for much larger SNR. The part of curves for $\mu < 2$ have some fluctuations due to the fact that the bias term $b$ varies a lot when SNR is small. When $\mu \geq 2$, the SVM has $b \approx 0$, and we can approximate the (strong-)adversarial-error rate by dropping the bias term in (8) and (12) and replace $\theta$ with its asymptotic limit as given by solving $(\theta, t)$ from the equations (Huang, 2017):

$$\sin^2 \theta = \frac{N}{d} \int_{-\infty}^{t} (t - z)^2 \varphi(z) \mathrm{d}z, \qquad \cos \theta = \frac{N}{d} \cdot \frac{\mu}{\sigma} \int_{-\infty}^{t} (t - z) \varphi(z) \mathrm{d}z \tag{13}$$

where $\varphi(z)$ is the density function of standard normal distribution. The rates plotted with these approximate formulas overlap the curves on Figure 2 very well for the part $\mu \geq 2$. We use these formulas to study the robustness of SVM against (strong-)adversarial examples.

Figure 3(a) plots the three error rates formulas of SVM when $\eta = 0.3$. Figure 3(b) plots the same rates for the Bayes classifier. These two classifiers are similar in misclassification rates and adversarial-error rates, but are very different in strong-adversarial-error rates. For a linear classier with small bias $b' \approx 0$, equations (6), (8) and (12) become:

$$p_m \approx 1 - \Phi\left(\frac{\|\boldsymbol{\mu}\|}{\sigma} \cos \theta\right), \quad p_{err} \approx 1 - \Phi\left[\left(\frac{\|\boldsymbol{\mu}\|}{\sigma} - \frac{\varepsilon}{\sigma}\right) \cos \theta\right], \quad p_{s-err} \approx 1 - \Phi\left[\left(\frac{\|\boldsymbol{\mu}\|}{\sigma} - \frac{\delta}{\sigma}\right) \cos \theta - \frac{\varepsilon}{\sigma} \sin \theta\right]$$
$$\tag{14}$$

Setting $\theta = 0$, we get the theoretical optimal rates achieved by the ideal Bayes classifier:

$$p_m^{id} = 1 - \Phi\left(\frac{\|\boldsymbol{\mu}\|}{\sigma}\right), \qquad p_{err}^{id} = 1 - \Phi\left(\frac{\|\boldsymbol{\mu}\|}{\sigma} - \frac{\varepsilon}{\sigma}\right), \qquad p_{s-err}^{id} = 1 - \Phi\left(\frac{\|\boldsymbol{\mu}\|}{\sigma} - \frac{\delta}{\sigma}\right). \tag{15}$$

Comparing equations (14) and (15), between the Bayes classifier and a linear classifier with small bias, both the misclassification rate and adversarial-error rate differ by a factor of $\cos \theta$ inside the $\Phi(\cdot)$ function. However, comparing
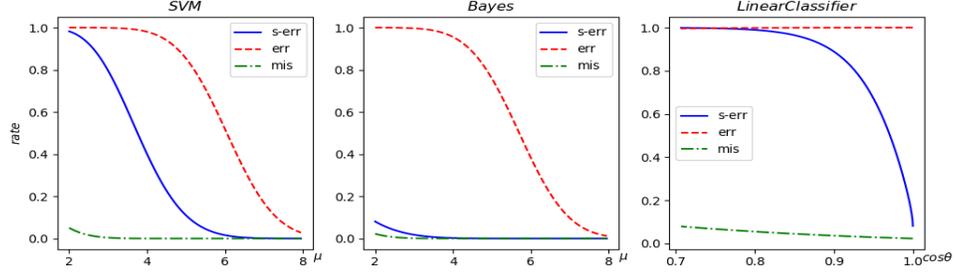
Figure 3: When $\eta = 0.3$, the three error rates (a) of SVM versus SNR $\mu$; (b) of Bayes classifier versus SNR $\mu$; (c) of an unbiased linear classifier versus $\cos\theta$ when $\mu = 2$.

their strong-adversarial-error rates, besides the multiplicative factor $\cos\theta$, there is also an extra bias term $-\frac{\varepsilon}{\sigma}\sin\theta$ inside the $\Phi(\cdot)$ function. Since $\frac{\varepsilon}{\sigma}$ is of order $O(\sqrt{d})$, $\theta = o(1/\sqrt{d})$ is needed for the linear classifier to approach the optimal strong-adversarial-error rate. In contrast, for the misclassification rate and adversarial-error rate, only $\theta = o(1)$ is needed to approach the optimal rates. Figure 3(c) plots these three rates versus $\cos\theta$ when $\eta = 0.3$ and $\mu = 2$. We can see that misclassification rate is low for a wide range of $\cos\theta$ values while the strong-adversarial-error rate only becomes low when $\cos\theta$ is very close to one.

## 3.2   Defending Against Strong-Adversarial Example Attacks

We have just seen that training a strong-adversarial-robust classifier needs stricter training requirements than those for a classifier with low misclassification rate: $\theta = o(1/\sqrt{d})$ versus $\theta = o(1)$. This is doable by incorporating some extra knowledge about the classification setting into the training. As an illustration, we show the results of using a naive method to find a sparse SVM in this case: for the SVM trained using standard method, takes ten non-zero components of $\boldsymbol{w}$ with largest absolute coefficients and set rest of components zero. The left panel of Figure 4 plots the strong-adversarial-error rates of this sparse SVM versus original SVM. We can see that the sparse SVM achieves a low strong-adversarial-error rate very close to the optimal rate of the ideal Bayes classifier. However, the same way of finding a sparse SVM does not produce strong-adversarial-robust classifier, shown in the right panel of Figure 4, when the data are generated with $\boldsymbol{\mu}_{+} = (\mu, \mu, ..., \mu)/\sqrt{d}$ instead of $\boldsymbol{\mu}_{+} = (\mu, 0, ..., 0)$. The data distributions in these two cases are equivalent with a change of coordinate systems. The sparse SVM fails in the second case since the extra knowledge incorporated into training is incorrect (sparseness only happens in the first coordinate system but not in the second coordinate system).
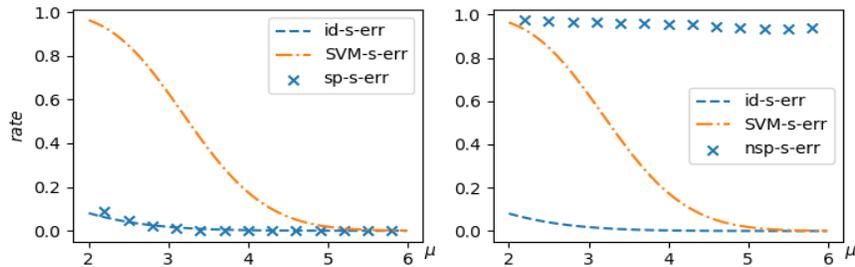


Figure 4: The strong-adversarial-error rates of standard SVM ($SVM - s - err$), the sparse SVM ($sp - s - err$) and the ideal Bayes classifier ($id - s - err$). Left: $\boldsymbol{\mu}_{+} = (\mu, 0, ..., 0)$; Right: $\boldsymbol{\mu}_{+} = (\mu, \mu, ..., \mu)/\sqrt{d}$. $\eta = 0.3$.

The above exercise shows that, even when adversarial examples are unavoidable, strong-adversarial-robust linear classifiers can be found with extra structural information on the underlying problem. Notice that the sparse SVM above provides good defense by using only the knowledge of a sparse representation existence (under the coordination system) but not what the sparse representation is, with the later part learned from data by training. More generally, statisticians have noticed that SVMs are suspect to the phenomenon of data-piling: there are more data points close to the decision

7

boundary than Gaussian mixture distribution implies. The distance-weighted discriminant (Marron et al., 2007) can be used to alleviate this data-piling phenomenon, and may be used to protect against strong-adversarial examples.
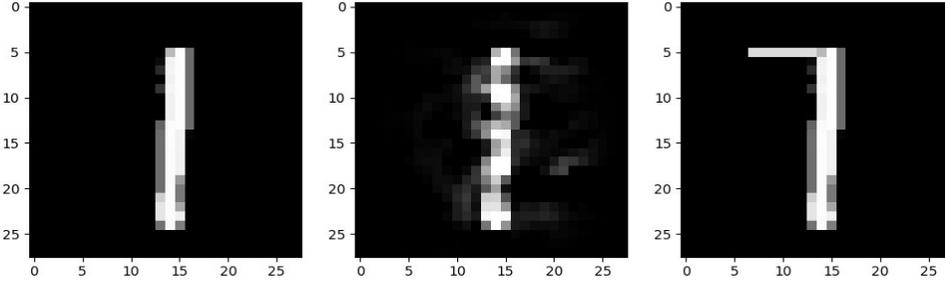


Figure 5: MNIST images of '1': (a) the original image, (b) an adversarial example with $\varepsilon = 2.34$, (c) a hand-made example with $\varepsilon = 2.28$

For more general classification problems, the signal direction is harder to define. But the concept of adversarial versus strong-adversarial examples still applies. Figure 5 shows an image of '1' from the MNIST data set, and two images with added perturbations. (b) shows an adversarial example obtained by Carlini and Wagner (2017) algorithm with $\varepsilon = 2.34$, that is misclassified by a DNN. (c) shows an image we made with a similar perturbation amount $\varepsilon = 2.29$. If a classifier is adversarial-robust at level of $\varepsilon = 2.34$, then it needs to classify both images (b) and (c) as '1'. However, classifying image (c) as '1' clearly contradicts what a human would do, rendering the usefulness of the classifier for practical applications in doubt. Generally, we should pursue a strong-adversarial-robust classifier, not an adversarial-robust one.

## 4    Discussions and Conclusions

In this paper, we provide clear definitions of adversarial and strong adversarial examples in the linear classification setting. Quantitative analysis shows that adversarial examples are hard to avoid but also should not be of concern in practice. Rather, we should focus on finding strong-adversarial-robust classifiers. We now consider the implications of these results on studying adversarial examples for general classifiers, and their relationship to some recent works in literature.

Recently, Shafahi et al. (2019) shows that no classifier can achieve low misclassification rate and also be adversarial-robust for data distributions with bounded density on a compact region in a high-dimensional space. Our analysis does not match exactly with their impossibility statement because we are studying the Gaussian mixture case, which has positive density on the whole space. However, in spirit our results have similar implications: for the usual SNR $O(1)$ that allows low misclassification rate, generally it is impossible to be also adversarial-robust (for which a much bigger SNR $O(\sqrt{d})$ is required).

Our results, however, do show that there can be adversarial-robust classifiers under the traditional definition when the SNR is very big. Schmidt et al. (2018) has also shown that, for Gaussian mixture classification problem and a particular training method, the adversarial-robustness is achievable but requires more training data than simply achieving the low misclassification rate only. Our formula indicates that useful adversarial-robust classifier do exist at the SNR level they assumed. Our study is more focused on the fundamental issue of when useful adversarial-robust classifiers exist, not which training method and what data complexity will find such a classifier. However, our formulas do indicate that an adversarial-robust classifier has to satisfy a stricter requirement than a good performing classifier. Thus either a better training method or a higher data complexity is needed for finding a useful adversarial-robust classifier, agreeing with the general theme of Schmidt et al. (2018).

Our results on the existence of adversarial examples do not change qualitatively when using other $\ell_p$ norm to measure the perturbation: under traditional definition, useful adversarial-robust classifier exists only when the data distribution has a very big SNR of $O(d^{min(1/p,1/2)})$ as shown in the Appendix 5. For many applications where good classifiers exists (SNR of only $O(1)$ ensures this), we can not pursue adversarial-robust classifier under the traditional adversarial example definition 1. The current defense strategies based on such adversarial example definition likely will still be suspect to more sophisticated adversarial attacks. For certifiable adversarial-robust classifiers (Madry et al., 2018; Sinha et al., 2018), the robustness is achieved only for the perturbation amount $\varepsilon$ high enough so that they differ from human in classifying images like those in Figure 1(c) and Figure 5(c). Thus a paradigm change is needed: we should train a classifier to be strong-adversarial-robust rather than adversarial-robust.

While the signal direction is obvious in the linear classification, the signal direction and the definition of strong-adversarial examples in general classification warrants further study. The signal direction in the linear classification here is the direction where the likelihood ratio of the two classes changes most rapidly. One reasonable extension is to define the signal direction at any data vector $x$ as the gradient direction of the likelihood ratio at $x$. Then similar to definition 2, the strong-adversarial example for general classifier also restrict the change along this signal direction to the amount $\delta$. The strong-adversarial-robust classifiers therefore are likely to be very close to the Bayes classifier. Some recent works have attempted training DNN to be close to the Bayes classifier: Wang et al. (2018) uses a nearest neighbors method, and Schott et al. (2019) applies the generative model techniques. In particular, Schott et al. (2019) applied their method on MNIST dataset, and when applying a specifically designed attack on such a trained DNN, the adversarial examples found are semantically meaningful for humans. That is, these adversarial examples are adversarial in traditional definition but likely not strong-adversarial. The new strong-adversarial examples framework can allow theoretical quantification of the robustness for these training methods. The analysis of strong-adversarial-robustness for general classifiers such as DNN can provide a new research direction on how to defend against realistic adversarial attacks.

## References

C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *CoRR*, vol. abs/1312.6199, 2014.

I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *International Conference on Learning Representations*, 2015. [Online]. Available: http://arxiv.org/abs/1412.6572

S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: A simple and accurate method to fool deep neural networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 IEEE Symposium on Security and Privacy (SP)*, May 2017, pp. 39–57.

A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *International Conference on Learning Representations*, 2018. [Online]. Available: https://openreview.net/forum?id=rJzIBfZAb

M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '16.   New York, NY, USA: ACM, 2016, pp. 1528–1540. [Online]. Available: http://doi.acm.org/10.1145/2976749.2978392

T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, "Adversarial patch," *arXiv preprint arXiv:1712.09665*, 2017.

A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *Artificial Intelligence Safety and Security*.   Chapman and Hall/CRC, 2018, pp. 99–112.

A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, "Synthesizing robust adversarial examples," in *International Conference on Machine Learning*, 2018, pp. 284–293.

N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," in *2016 IEEE Symposium on Security and Privacy (SP)*.   IEEE, 2016, pp. 582–597.

A. Sinha, H. Namkoong, and J. Duchi, "Certifiable distributional robustness with principled adversarial training," in *International Conference on Learning Representations*, 2018. [Online]. Available: https://openreview.net/forum?id=Hk6kPgZA-

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

A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," in *International Conference on Machine Learning*, 2018, pp. 274–283.

A. Shafahi, W. R. Huang, C. Studer, S. Feizi, and T. Goldstein, "Are adversarial examples inevitable?" in *International Conference on Learning Representations*, 2019. [Online]. Available: https://openreview.net/forum?id=r1lWUoA9FQ

A. Fawzi, O. Fawzi, and P. Frossard, "Analysis of classifiers' robustness to adversarial perturbations," *Machine Learning*, vol. 107, no. 3, pp. 481–508, Mar 2018. [Online]. Available: https://doi.org/10.1007/s10994-017-5663-3

H. Huang, "Asymptotic behavior of support vector machine for spiked population model," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 1472–1492, 2017.

J. S. Marron, M. J. Todd, and J. Ahn, "Distance-weighted discrimination," *Journal of the American Statistical Association*, vol. 102, no. 480, pp. 1267–1271, 2007. [Online]. Available: http://www.jstor.org/stable/27639976

L. Schmidt, S. Santurkar, D. Tsipras, K. Talwar, and A. Madry, "Adversarially robust generalization requires more data," in *Advances in Neural Information Processing Systems*, 2018, pp. 5014–5026.

Y. Wang, S. Jha, and K. Chaudhuri, "Analyzing the robustness of nearest neighbors to adversarial examples," in *International Conference on Machine Learning*, 2018, pp. 5120–5129.

L. Schott, J. Rauber, M. Bethge, and W. Brendel, "Towards the first adversarially robust neural network model on mnist," in *Seventh International Conference on Learning Representations (ICLR 2019)*, 2019, pp. 1–16.

G. O'Brien, "Limit theorems for the maximum term of a stationary process," *The Annals of Probability*, pp. 540–545, 1974.

## 5   Appendix

### 5.1   Proof of Lemma 1

**Lemma 2.** *The defining sets for $\varepsilon$-adversarial and $(\varepsilon, \delta)$-strong adversarial examples are given by:*

$$\Omega_\varepsilon = \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0); \quad \Omega_{\varepsilon,\delta} = \Omega(\boldsymbol{u}_2) \cup \Omega(-\boldsymbol{u}_2) \tag{16}$$

*where $\boldsymbol{u}_2 = \beta\boldsymbol{\mu}_0 + \sqrt{\varepsilon^2 - \beta^2}\boldsymbol{n}_0, \beta = \min(\varepsilon\cos\theta, \delta)$.*

*Proof.* **Proof of the adversarial defining set formula.**   Since it is obvious from the definition that $\Omega_\varepsilon = \bigcup_{\|\boldsymbol{v}\|\leq\varepsilon} \Omega(\boldsymbol{v}) \supseteq \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0)$, we only need to show that $\Omega_\varepsilon \subseteq \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0)$. That is, for any data point $\boldsymbol{x} \in \Omega_\varepsilon$, either $\boldsymbol{x} + \varepsilon\boldsymbol{v}_0$ or $\boldsymbol{x} - \varepsilon\boldsymbol{v}_0$ changes its classification.

We now claim that the last statement is equivalent to that $\varepsilon\boldsymbol{v}_0$ is the solution to the optimization problem:

$$\max \boldsymbol{w} \cdot \boldsymbol{v}, \qquad \boldsymbol{v} \in D_\varepsilon = \{\boldsymbol{v} \in \mathbb{R}^d : \|\boldsymbol{v}\| \leq \varepsilon\}. \tag{17}$$

To see this, if $\varepsilon\boldsymbol{v}_0$ is the solution, then $\boldsymbol{w} \cdot \boldsymbol{v} \leq \boldsymbol{w} \cdot (\varepsilon\boldsymbol{v}_0) = \varepsilon\|\boldsymbol{w}\|$ for all $\boldsymbol{v} \in D_1$. Now for a $\boldsymbol{x}$ classified into the '$-$' class and $\boldsymbol{x} \in \Omega_\varepsilon$, then $\boldsymbol{w} \cdot \boldsymbol{x} + b < 0$ and $\boldsymbol{w} \cdot (\boldsymbol{x} + \boldsymbol{v}) + b > 0$. Hence

$$\boldsymbol{w} \cdot (\boldsymbol{x} + \varepsilon\boldsymbol{v}_0) + b \geq \boldsymbol{w} \cdot \boldsymbol{x} + \boldsymbol{w} \cdot \boldsymbol{v} + b > 0,$$

that is, $\boldsymbol{x} + \varepsilon\boldsymbol{v}_0$ is misclassified into the '$+$' class thus $\boldsymbol{x} \in \Omega(\varepsilon\boldsymbol{v}_0)$. By symmetry, $-\varepsilon\boldsymbol{v}_0$ is the solution to $\min \boldsymbol{w} \cdot \boldsymbol{v}$ when $\boldsymbol{v} \in D_1$, and hence $-\varepsilon\|\boldsymbol{w}\| \leq \boldsymbol{w} \cdot \boldsymbol{v}$ also for all $\boldsymbol{v} \in D_1$. Hence for a $\boldsymbol{x}$ classified into the '$+$' class and $\boldsymbol{x} \in \Omega_\varepsilon$, similarly we have that $\boldsymbol{x} - \varepsilon\boldsymbol{v}_0$ is misclassified into the '$-$' class thus $\boldsymbol{x} \in \Omega(-\varepsilon\boldsymbol{v}_0)$.

Finally, $\varepsilon\boldsymbol{v}_0$ is indeed the solution to (17) due to the Cauchy-Schwartz inequality $\boldsymbol{w} \cdot \boldsymbol{v} \leq \|\boldsymbol{w}\|\|\boldsymbol{v}\| \leq \|\boldsymbol{w}\|\varepsilon$. The first equality holds if and only if $\boldsymbol{v}$ is along the same direction of $\boldsymbol{w}$, thus $\boldsymbol{v} = c\boldsymbol{v}_0$. The second equality holds if and only if $\|\boldsymbol{v}\| = \varepsilon$, thus $\boldsymbol{v} = \varepsilon\boldsymbol{v}_0$. This finishes the proof for $\Omega_\varepsilon = \Omega(\varepsilon\boldsymbol{v}_0) \cup \Omega(-\varepsilon\boldsymbol{v}_0)$.

**Proof of the strong adversarial defining set formula.** The proof follows exactly the outline of the adversarial case proof above. Only now we need to prove that $\boldsymbol{u}_2$ is the solution to the optimization problem

$$\max \boldsymbol{w} \cdot \boldsymbol{v}, \qquad \boldsymbol{v} \in D_{\varepsilon,\delta} = \{\boldsymbol{v} \in \mathbb{R}^d : \|\boldsymbol{v}\| \leq \varepsilon, |\boldsymbol{v} \cdot \boldsymbol{\mu}_0| \leq \delta\}. \tag{18}$$

We can decompose $\boldsymbol{w}$ as $\boldsymbol{w} = (\boldsymbol{w} \cdot \boldsymbol{\mu}_0)\boldsymbol{\mu}_0 + (\boldsymbol{w} \cdot \boldsymbol{n}_0)\boldsymbol{n}_0$, accordingly, $\boldsymbol{v}$ can be decomposed as $\boldsymbol{v} = (\boldsymbol{v} \cdot \boldsymbol{\mu}_0)\boldsymbol{\mu}_0 + (\boldsymbol{v} \cdot \boldsymbol{n}_0)\boldsymbol{n}_0 + (\boldsymbol{v} \cdot \boldsymbol{m}_0)\boldsymbol{m}_0$, where $\boldsymbol{m}_0$ is the unit normal vector of the plane spanned by $\boldsymbol{\mu}_0$ and $\boldsymbol{w}$, therefore

$$\boldsymbol{w} \cdot \boldsymbol{v} = (\boldsymbol{w} \cdot \boldsymbol{\mu}_0)(\boldsymbol{v} \cdot \boldsymbol{\mu}_0) + (\boldsymbol{w} \cdot \boldsymbol{n}_0)(\boldsymbol{v} \cdot \boldsymbol{n}_0) = \cos\theta(\boldsymbol{v} \cdot \boldsymbol{\mu}_0) + \sin\theta(\boldsymbol{v} \cdot \boldsymbol{n}_0) := x\cos\theta + y\sin\theta. \tag{19}$$

The optimization problems becomes to maximize $x\cos\theta + y\sin\theta$ in (19) under the constraints $x^2 + y^2 = \varepsilon^2 - (\boldsymbol{v} \cdot \boldsymbol{m}_0)^2, |x| \leq \delta$. This is a linear programming setup, it is easy to see that first we must have $\boldsymbol{v} \cdot \boldsymbol{m}_0 = 0$ to reach maximum. Then the solution is either at the corner $(x, y) = (\delta, \sqrt{\varepsilon^2 - \delta^2})$ or at the tangent point $(x, y) = \varepsilon(\cos\theta, \sin\theta)$ as in semi-adversarial case. If $\varepsilon\cos\theta < \delta$, the solution is at the tangent point $(x, y) = \varepsilon(\cos\theta, \sin\theta)$. Otherwise, the solution is at the corner $(x, y) = (\delta, \sqrt{\varepsilon^2 - \delta^2})$. Combining the two cases, we arrive at the formula for $\boldsymbol{u}_2$ under equation (16).

$\square$

### 5.2   $\ell_p$-Adversarial and $\ell_p$-Strong-Adversarial Rates

In literature, the adversarial examples have been studied under different norms. Here we extend the analysis in main text to the general $\ell_p$ norms with $p \in [1, \infty]$[3]. That is, we use the distance metric $d_p(\boldsymbol{x}, \boldsymbol{y}) = \|\boldsymbol{x} - \boldsymbol{y}\|_p$. Also, we denote $\ell_q$ as the dual of $\ell_p$, i.e., $1/p + 1/q = 1$.

Therefore the classical adversarial examples definition becomes the following.

**Definition 3.** *Given a classifier $C$, an $\varepsilon$-$\ell_p$-adversarial example of a data vector $\boldsymbol{x}$ is another data vector $\boldsymbol{x}'$ such that $d_p(\boldsymbol{x}, \boldsymbol{x}') \le \varepsilon$ but $C(\boldsymbol{x}) \neq C(\boldsymbol{x}')$.*

As before, we restrict the perturbation amount along the signal direction $\boldsymbol{\mu}_0$ to $\delta$ for strong-adversarial examples.

**Definition 4.** *Given a classifier $C$, an $(\varepsilon, \delta)$-$\ell_p$-strong-adversarial example of a data vector $\boldsymbol{x}$ is another data vector $\boldsymbol{x}'$ such that $d_p(\boldsymbol{x}, \boldsymbol{x}') \le \varepsilon$ and $|(\boldsymbol{x} - \boldsymbol{x}') \cdot \boldsymbol{\mu}_0| \le \delta$ but $C(\boldsymbol{x}) \neq C(\boldsymbol{x}')$.*

### 5.3   $\ell_p$-Adversarial Rate and Existence of $\ell_p$-Adversarial-Robust Classifiers

The analysis follows the same outline as the analysis for the $\ell_2$ norm case. We first characterize the defining set $\Omega_{\varepsilon|p} = \{\boldsymbol{x} : \boldsymbol{x} \text{ has an } \varepsilon - \ell_p\text{-adversarial example}\}$.

**Lemma 3.** *The defining sets for $\varepsilon$-$\ell_p$-adversarial examples is given by:*

$$\Omega_{\varepsilon|p} = \Omega(\varepsilon \boldsymbol{v}_{0|p}) \cup \Omega(-\varepsilon \boldsymbol{v}_{0|p}) \tag{20}$$

*where $\boldsymbol{v}_{0|p}$ is the $d$-dimensional vector with component $(\boldsymbol{v}_{0|p})_i = \mathrm{sgn}(w_i) \cdot (|w_i| / \|\boldsymbol{w}\|_q)^{q-1}$.*

Here sgn denotes the sign function. That is, $\mathrm{sgn}(x) = 1$ for $x > 0$; $\mathrm{sgn}(x) = -1$ for $x < 0$ and $\mathrm{sgn}(0) = 0$.

Furthermore, we denote the $p$-th power of a vector $\boldsymbol{v} = (v_1, ..., v_d)$ as taking the power component-wise. That is, $(\boldsymbol{v}^p)_i = \mathrm{sgn}(v_i) \cdot |v_i|^p$. Then the above $\boldsymbol{v}_{0|p}$ can be rewritten as $\boldsymbol{v}_{0|p} = (\boldsymbol{w} / \|\boldsymbol{w}\|_q)^{q-1}$.

*Proof.* The proof is similar to the proof of Lemma 2. Following the derivations there, we only need to show that $\boldsymbol{v}_0 = (\boldsymbol{w} / \|\boldsymbol{w}\|_q)^{q-1}$ is the solution to the optimization problem:

$$\max |\boldsymbol{w} \cdot \boldsymbol{v}|, \qquad \boldsymbol{v} \in D_{\varepsilon|p} = \{\boldsymbol{v} \in \mathbb{R}^d : \|\boldsymbol{v}\|_p \le \varepsilon\}. \tag{21}$$

By Holder's inequality, we have $|\boldsymbol{w} \cdot \boldsymbol{v}| \le \|\boldsymbol{w}\|_q \|\boldsymbol{v}\|_p \le \varepsilon \|\boldsymbol{w}\|_q$.

For the first "$\le$" to be "$=$", $\boldsymbol{v}^p$ has to be proportional to $\boldsymbol{w}^q$. That is, for some constant $c$, $\boldsymbol{v} = c\boldsymbol{w}^{q/p} = c\boldsymbol{w}^{q-1}$. For the second "$\le$" to be "$=$", we need $\varepsilon = \|\boldsymbol{v}\|_p$. That is,

$$\varepsilon^p = \|\boldsymbol{v}\|_p^p = c^p \sum_{i=1}^d |v_i|^p = c^p \sum_{i=1}^d (|w_i|^{q/p})^p = c^p \sum_{i=1}^d (|w_i|^q) = c^p \|\boldsymbol{w}\|_q^q.$$

Hence we have $\varepsilon = c \|\boldsymbol{w}\|_q^{q/p} = c \|\boldsymbol{w}\|_q^{q-1}$, and thus $c = \varepsilon \|\boldsymbol{w}\|_q^{1-q}$. Plug $c$ into $\boldsymbol{v} = c\boldsymbol{w}^{q-1}$, we get $\boldsymbol{v} = \varepsilon(\boldsymbol{w} / \|\boldsymbol{w}\|_q)^{q-1} = \varepsilon \boldsymbol{v}_{0|p}$. This is the solution to the optimization problem. Hence arguments similar to those for the proof of Lemma 2 above show that the equation (20) gives the defining set here.  $\square$

With the characterization lemma 3, we can then compute the adversarial rate as before. Note that the misclassification rate has nothing to do with the perturbation for adversarial examples. Thus regardless of which $\ell_p$ norm is used to measure the perturbation, the misclassification rate is still given by the same formula as before.

$$p_m = 1 - 0.5 \left[ \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} + b'}{\|\boldsymbol{w}\| \sigma} \right) + \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} - b'}{\|\boldsymbol{w}\| \sigma} \right) \right]. \tag{22}$$

The calculation of $\ell_p$-adversarial rate follows $\ell_2$-adversarial rate calculation exactly, except that the term $\varepsilon \|\boldsymbol{w}\|_2$ is replaced by $\boldsymbol{w} \cdot \varepsilon \boldsymbol{v}_{0|p} = \varepsilon \|\boldsymbol{w}\|_q$. Therefore, we have the following result.

**Theorem 3.** *The overall $\ell_p$-adversarial rate of a linear classifier for the balanced Gaussian mixture data is*

$$p_{adv|p} = 1 - p_m - 0.5 \left[ \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} + b'}{\|\boldsymbol{w}\|_2 \sigma} - \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2} \frac{\varepsilon}{\sigma} \right) + \Phi\left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} - b'}{\|\boldsymbol{w}\|_2 \sigma} - \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2} \frac{\varepsilon}{\sigma} \right) \right]. \tag{23}$$

---

[3]We did not consider $p \in [0, 1)$ because in this case, $d_p$ is not a metric, although practically, $\ell_0$ is considered.

Now we have the $\ell_p$-adversarial error formula as the following.

$$
\begin{aligned}
p_{err|p} = p_{adv|p} + p_m &= 1 - 0.5\left[\Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}+b'}{\|\boldsymbol{w}\|_2\,\sigma} - \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}\frac{\varepsilon}{\sigma}\right) + \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}-b'}{\|\boldsymbol{w}\|_2\,\sigma} - \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}\frac{\varepsilon}{\sigma}\right)\right] \\
&\geq 1 - \Phi\left(\frac{\boldsymbol{w}\cdot\boldsymbol{\mu}}{\|\boldsymbol{w}\|_2\,\sigma} - \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}\frac{\varepsilon}{\sigma}\right) = 1 - \Phi\left(\frac{\|\boldsymbol{\mu}\|_2}{\sigma}\cos\theta - \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}\frac{\varepsilon}{\sigma}\right)
\end{aligned}
\tag{24}
$$

Corresponding to the discussions in the main text, a useful classifier only requires a signal-noise ratio (SNR) of $\|\mu\|_2/\sigma = O(1)$ due to equation (22).

In contrast, due to equation (24), a necessary condition for the existence of a $\ell_p$-adversarial-robust classifier is

$$
\frac{\|\boldsymbol{\mu}\|_2}{\sigma} - \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}\frac{\varepsilon}{\sigma} = O(1).
\tag{25}
$$

We now investigate what order of SNR $\|\mu\|_2/\sigma$ is needed to make (25) hold.

First, we have to find the practical relevant order of $\varepsilon$ needs to be studied. The following lemma about the average $\ell_p$-norm of Gaussian noise will provide the guideline.

**Lemma 4.** *Let the random vector $\boldsymbol{x} = (x_1, ..., x_d)$ follows the d-dimensional Gaussian distribution $N(0, \sigma^2 I_d)$. Then*

$$
\begin{cases}
E[\|x\|_p^p] = m_p d\sigma^p, & p \in [1, \infty), \\
E[\|x\|_p] = O(\sqrt{\log d}\,\sigma), & p = \infty,
\end{cases}
\tag{26}
$$

*where $m_p$ denotes the p-th moment of the standard Gaussian distribution.*

*Proof.* For $p \in [1, \infty)$,

$$
E[\|x\|_p^p] = E[\sum_{i=1}^d |x_i|^p] = \sum_{i=1}^d E[|x_i|^p] = dE[|x_1|^p] = d\sigma^p m_p.
\tag{27}
$$

The $\ell_\infty$ result follows directly from the large deviation formula obtained by O'Brien (1974). □

In the Gaussian mixture data, Lemma 4 states that the average $\ell_p$ noise is $d^{1/p}\sigma m_p^{1/p}$. Therefore, for an $\eta < 1$, a perturbation amount of $\varepsilon = \eta d^{1/p}\sigma m_p^{1/p}$ will be smaller than the average noise thus hard to distinguish from noise (unless it concentrates in the signal direction). Thus any practical relevant defense needs to be robust at the minimum against perturbations of order $\varepsilon = O(d^{1/p}\sigma)$. For the $\ell_\infty$, the defense needs to be robust at the minimum against perturbations of order $\varepsilon = O(\sqrt{\log d}\,\sigma)$.

Plug-in these $\varepsilon$ orders into the necessary condition (25), the existence of a $\ell_p$-adversarial-robust classifier requires at least $\frac{\|\boldsymbol{\mu}\|_2}{\sigma} = O(\frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}d^{1/p})$ for $p \in [1, \infty)$; and it requires at least $\frac{\|\boldsymbol{\mu}\|_2}{\sigma} = O(\frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}\sqrt{\log d})$ for $p = \infty$.

Next we use the norm comparison inequality to find these orders. For any $\boldsymbol{w} \in \mathbb{R}^d$ and any $0 < r < s < \infty$, we have

$$
\|\boldsymbol{w}\|_s \leq \|\boldsymbol{w}\|_r \leq d^{1/r-1/s}\|\boldsymbol{w}\|_s.
\tag{28}
$$

**(A)** For $1 \leq p < 2$, we have $2 < q \leq \infty$. Using $r = 2$ and $s = q$ in (28), we get

$$
d^{1/q-1/2} \leq \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2} \leq 1.
$$

Plug this lower bound into the required order $\frac{\|\boldsymbol{\mu}\|_2}{\sigma} = O(\frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2}d^{1/p})$, the existence of a $\ell_p$-adversarial-robust classifier requires SNR of at least

$$
\frac{\|\boldsymbol{\mu}\|_2}{\sigma} = O(d^{1/q-1/2}d^{1/p}) = O(d^{1/p+1/q-1/2}) = O(d^{1/2}).
$$

**(B)** For $2 < p < \infty$, then $q < 2 < \infty$. Using $r = q$ and $s = 2$ in (28), we get

$$
1 \leq \frac{\|\boldsymbol{w}\|_q}{\|\boldsymbol{w}\|_2} \leq d^{1/q-1/2}.
$$

Thus the existence of a $\ell_p$-adversarial-robust classifier requires SNR of at least

$$\frac{\|\boldsymbol{\mu}\|_2}{\sigma} = O(1 \cdot d^{1/p}) = O(d^{1/p}).$$

**(C)** When $p = \infty$, then $q = 1$. Using $r = q$ and $s = 2$ in (28), we get

$$1 \le \frac{\|\boldsymbol{w}\|_1}{\|\boldsymbol{w}\|_2} \le d^{1/q - 1/2}.$$

Thus the existence of a $\ell_p$-adversarial-robust classifier requires SNR of at least

$$\frac{\|\boldsymbol{\mu}\|_2}{\sigma} = O(1 \cdot \sqrt{\log d}) = O(\sqrt{\log d}).$$

We summarize the results for cases (A), (B) and (C) into the following theorem.

**Theorem 4.** *For linear classification of balanced Gaussian mixture data, the existence of a $\ell_p$-adversarial-robust classifier requires SNR of at least*

$$\frac{\|\boldsymbol{\mu}\|_2}{\sigma} = \begin{cases} O(d^{\min(1/p, 1/2)}) & p \in [1, \infty), \\ O(\sqrt{\log d}) & p = \infty. \end{cases} \tag{29}$$

Theorem 4 shows that the required SNR magnitude for $\ell_p$-adversarial-robustness differs for different $p$. The $\ell_p$-adversarial-robustness is hardest to achieve for $1 \le p \le 2$ since the required SNR $O(\sqrt{d})$ is highest in these cases. The $\ell_\infty$-adversarial-robustness has the smallest SNR requirement, thus easiest to achieve. This agrees with the observation by Schott et al. (2019): the $\ell_\infty$ robust classifier in Madry et al. (2018) is still highly susceptible to $\ell_2$ attack.

### 5.4 $\ell_p$-Strong-Adversarial Rate and Existence of $\ell_p$-Strong-Adversarial-Robust Classifiers

Following the same derivations before, we have the following lemma for the defining set $\Omega_{\varepsilon, \delta|p} = \{\boldsymbol{x} : \boldsymbol{x} \text{ has an } (\varepsilon, \delta) - \ell_p\text{-strong-adversarial example}\}$.

**Lemma 5.** *The defining set for $(\varepsilon, \delta) - \ell_p$-strong-adversarial examples is given by:*

$$\Omega_{\varepsilon, \delta|p} = \Omega(\boldsymbol{u}_p) \cup \Omega(-\boldsymbol{u}_p) \tag{30}$$

*where $\boldsymbol{u}_p$ is the solution to the optimization problem:*

$$\max |\boldsymbol{w} \cdot \boldsymbol{v}|, \qquad \boldsymbol{v} \in D_{\varepsilon, \delta|p} = \{\boldsymbol{v} \in \mathbb{R}^d : \|\boldsymbol{v}\|_p \le \varepsilon, |\boldsymbol{v} \cdot \boldsymbol{\mu}_0| \le \delta\} \subset D_{\varepsilon|p}. \tag{31}$$

Notice the optimization problem of $\max |\boldsymbol{w} \cdot \boldsymbol{v}|$ is a linear programming problem, and the feasible region $D_{\varepsilon, \delta|p}$ is a convex region. Therefore the solution $\boldsymbol{u}_p$ does exist.

Now replace the term $\varepsilon \|\boldsymbol{w}\|_2$ by $\boldsymbol{w} \cdot \boldsymbol{u}_p$ in the previous derivations of $(\varepsilon, \delta)$-strong-adversarial rate using $\ell_2$ norm, we get the following Theorem.

**Theorem 5.** *The overall $(\varepsilon, \delta) - \ell_p$-strong-adversarial rate of a linear classifier for the balanced Gaussian mixture data is*

$$p_{s-adv|p} = 1 - p_m - 0.5 \left[ \Phi \left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} + b'}{\|\boldsymbol{w}\|_2 \sigma} - \frac{\boldsymbol{w} \cdot \boldsymbol{u}_p}{\|\boldsymbol{w}\|_2 \sigma} \right) + \Phi \left( \frac{\boldsymbol{w} \cdot \boldsymbol{\mu} - b'}{\|\boldsymbol{w}\|_2 \sigma} - \frac{\boldsymbol{w} \cdot \boldsymbol{u}_p}{\|\boldsymbol{w}\|_2 \sigma} \right) \right] \tag{32}$$

We now try to find an SNR order that allows $\ell_p$-strong-adversarial-robustness by applying formula (32) to the Bayes classifier whose $\boldsymbol{w} = \boldsymbol{\mu}_0$ and $b' = 0$. In this case, the solution to the optimization problem (31) becomes $\boldsymbol{u}_p = \delta \boldsymbol{\mu}_0$. Thus using formula (32) we can get the $(\varepsilon, \delta) - \ell_p$-strong-adversarial-error rate for the Bayes classifier as

$$p_{s-err|p} = 1 - \Phi \left( \frac{\|\boldsymbol{\mu}\|_2 - \delta}{\sigma} \right). \tag{33}$$

Since practical relevant $\delta$ can not exceed $\sigma$ (in that case, no classifier can work as at least half of all data vectors will be perturbed into another class), SNR $\frac{\|\boldsymbol{\mu}\|_2}{\sigma}$ of order $O(1)$ can result in a useful $\ell_p$-strong-adversarial-robust classifier. This agrees with the conclusions in the main text about the existence of $\ell_2$-strong-adversarial-robust classifiers.