
On the Diversity of Adversarial Ensemble Learning

Jun-Qi Guo^{*1} Meng-Zhang Qian^{*1} Wei Gao¹ Zhi-Hua Zhou¹

Abstract

Diversity has been one of the most crucial factors on the design of adversarial ensemble methods. This work focuses on the fundamental problems: How to define the diversity for the adversarial ensemble, and how to correlate with algorithmic performance. We first show that it is an NP-Hard problem to precisely calculate the diversity of two networks in adversarial ensemble learning, which makes it different from prior diversity analysis. We present the first diversity decomposition under the first-order approximation for the adversarial ensemble learning. Specifically, the adversarial ensemble loss can be decomposed into average of individual adversarial losses, *gradient diversity*, *prediction diversity* and *cross diversity*. Hence, it is not sufficient to merely consider the gradient diversity on the characterization of diversity as in previous adversarial ensemble methods. We present diversity decomposition for classification with cross-entropy loss similarly. Based on the theoretical analysis, we develop new ensemble method via orthogonal adversarial predictions to simultaneously improve gradient diversity and cross diversity. We finally conduct experiments to validate the effectiveness of our method.

1. Introduction

General machine learning models may be misled heavily by examples with adversarial perturbations (Szegedy et al., 2014; Goodfellow et al., 2015), which raises some serious concerns about reliability of models, particularly in high-risk applications such as healthcare, finance and autonomous driving (Finlayson et al., 2019; Deng et al., 2021; Fursov et al., 2021). Various robust methods have been developed against adversarial examples in recent years (Zheng et al.,

2019b; Goyal et al., 2021; Pang et al., 2021; Wang et al., 2023; Peng et al., 2023; Bartoldson et al., 2024).

Ensemble learning combines multiple learners rather than one single learner with better performance, which has been paid much attention in the adversarial robustness learning. Sequential robust ensemble has been constructed via the boosting framework (Abernethy et al., 2021; Zhang et al., 2019a; 2022; Guo et al., 2022), and parallel ensemble has also developed for robust learning (Pinot et al., 2020; Sen et al., 2020; Yang et al., 2021; 2022; Deng & Mu, 2024).

Diversity has always been one of the most crucial factors in the design of ensemble methods (Zhou, 2012; Wood et al., 2023). For adversarial ensemble learning, Pang et al. (2019) took prediction disagreements of base learners as diversity, while Yang et al. (2020) defined the diversity via losses of base learners on exchanged adversarial examples. Another idea is to consider the misalignment of gradient directions for diversity (Kariyappa & Qureshi, 2019; Dabouei et al., 2020; Huang et al., 2021; Bogun et al., 2022).

There are some fundamental problems open for adversarial ensemble learning. For example, how to formally define the diversity of adversarial ensemble, and what's more, how to correlate diversity definition with algorithmic performance from a theoretical view. This work studies fundamental problems of diversity in the adversarial ensemble learning, and the main contributions are summarized as follows:

- We first show that it is an NP-Hard problem to precisely calculate the diversity of two neural networks in the adversarial ensemble, since diversity is heavily relevant to intrinsic structures and output predictions of models simultaneously. This challenge makes it different from traditional diversity on output predictions (Zhou, 2012; Wood et al., 2023). Sun & Zhou (2018) indicated the importance of structure diversity for decision trees, whereas it remains open for neural networks.
- We present the first diversity decomposition under the first-order approximation in the adversarial ensemble learning. Specifically, the adversarial ensemble loss is decomposed into average of individual adversarial losses, *prediction diversity*, *gradient diversity* and *cross diversity*. It is not sufficient to only consider gradient diversity on the characterization of diversity as in prior

^{*}Equal contribution ¹National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China; School of Artificial Intelligence, Nanjing University, Nanjing, China. Correspondence to: Wei Gao <gaow@lamda.nju.edu.cn>.

adversarial ensemble methods. We present similar decomposition for classification with cross-entropy loss, which is commonly used for neural networks.

- Based on theoretical analysis, we develop the AdvE_{OAP} adversarial ensemble method¹ via the orthogonal of adversarial predictions of base learners, which could improve gradient and cross diversity simultaneously. We finally conduct empirical studies to validate the effectiveness of our AdvE_{OAP} method.

The rest of this work is constructed as follows: Section 2 presents some preliminaries. Section 3 provides diversity analysis. Section 4 develops our method. Section 5 conducts experiments. Section 6 concludes with future work.

2. Preliminary

Let $\mathcal{X} \subseteq \mathbb{R}^d$ and \mathcal{Y} denote the instance and label space, respectively, where $\mathcal{Y} = \{0, 1\}$ for binary classification and $\mathcal{Y} \subseteq \mathbb{R}$ for regression. Let \mathcal{D} be an underlying distribution over the product space $\mathcal{X} \times \mathcal{Y}$, and we have a training data

$$S_n = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\},$$

where each sample is drawn i.i.d. from distribution \mathcal{D} .

Define the perturbation set Δ_p^ϵ w.r.t. l_p norm and $\epsilon > 0$ as

$$\Delta_p^\epsilon = \{\boldsymbol{\delta} \in \mathbb{R}^d : \|\boldsymbol{\delta}\|_p = (|\delta_1|^p + |\delta_2|^p + \dots + |\delta_d|^p)^{1/p} \leq \epsilon\},$$

which shows instances' imperceptible perturbation (Szegedy et al., 2014). Let $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$ be a hypothesis space, and loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is introduced to measure performance. Given $f \in \mathcal{F}$, we define the adversarial perturbation w.r.t. example \mathbf{x} as

$$\boldsymbol{\delta}_f^* \in \arg \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \ell(f(\mathbf{x} + \boldsymbol{\delta}), y),$$

and $\mathbf{x} + \boldsymbol{\delta}_f^*$ is called the adversarial example w.r.t. \mathbf{x} .

We define the expected adversarial loss as

$$\mathcal{L}^{\text{adv}}(f, \mathcal{D}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \{\ell(f(\mathbf{x} + \boldsymbol{\delta}), y)\} \right],$$

and define the empirical adversarial loss w.r.t. data S_n as

$$\hat{\mathcal{L}}^{\text{adv}}(f, S_n) = \frac{1}{n} \sum_{i=1}^n \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \{\ell(f(\mathbf{x}_i + \boldsymbol{\delta}), y_i)\}.$$

We finally introduce some useful notations in this work. Denote by $\langle \cdot, \cdot \rangle$ the inner product of two vectors, and \mathbf{e}_i is a unit vector with i -th element 1. Write $[k] = \{1, 2, \dots, k\}$

¹Code is available at <https://github.com/GuoJQ42/AdvOAP>.

for integer $k > 0$. For two non-negative real numbers a and b with $a + b = 1$, we define the KL-divergence as

$$\text{KL}(a, b) = a \ln(a/b) + (1 - a) \ln((1 - a)/(1 - b)),$$

and for two probability vectors $\mathbf{a} = (a_1, a_2, \dots, a_m)$ and $\mathbf{b} = (b_1, b_2, \dots, b_m)$, we define the KL-divergence as

$$\text{KL}(\mathbf{a}, \mathbf{b}) = \sum_{i=1}^m a_i \ln(a_i/b_i).$$

3. Theoretical Analysis on Diversity

Given m learners $f_1(\mathbf{x}), \dots, f_m(\mathbf{x})$, this work focuses on the simplest ensemble method (Dietterich, 2000; Zhou, 2012)

$$\bar{f}(\mathbf{x}) = \sum_{j=1}^m \frac{f_j(\mathbf{x})}{m}.$$

We begin with the challenge on the analysis of adversarial diversity, and then present diversity decomposition w.r.t. squared loss and cross-entropy loss, respectively.

3.1. Main challenge on analysis of adversarial diversity

For traditional non-adversarial ensemble learning, it is easy to make the error-ambiguity decomposition over example (\mathbf{x}, y) w.r.t. squared loss from (Krogh & Vedelsby, 1994; Zhou, 2012) as follows:

$$\underbrace{(\bar{f}(\mathbf{x}) - y)^2}_{\text{ensemble loss}} = \underbrace{\sum_{j=1}^m \frac{(f_j(\mathbf{x}) - y)^2}{m}}_{\text{average loss}} - \underbrace{\sum_{j=1}^m \frac{(f_j(\mathbf{x}) - \bar{f}(\mathbf{x}))^2}{m}}_{\text{ambiguity}},$$

where the ambiguity can be viewed as ensemble diversity. Wood et al. (2023) further presented bias-variance-diversity decomposition for ensemble learning, simplified by

$$\text{ensemble loss} = \text{bias} + \text{variance} - \text{diversity}.$$

It is natural to consider some similar decompositions in the adversarial diversity learning. However, this remains some challenges as shown by the following theorem.

Theorem 3.1. *For squared loss, it is an NP-hard problem to precisely calculate the diversity w.r.t. example (\mathbf{x}, y) in the adversarial ensemble learning as follows:*

$$\frac{1}{2} \sum_{j=1}^2 \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \{(f_j(\mathbf{x} + \boldsymbol{\delta}) - y)^2\} - \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \{(\bar{f}(\mathbf{x} + \boldsymbol{\delta}) - y)^2\},$$

where the error ambiguity decomposition is considered for two neural networks $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ with ReLU activation, and $\bar{f}(\mathbf{x}) = (f_1(\mathbf{x}) + f_2(\mathbf{x}))/2$.

Theorem 3.1 shows that, even for the ensemble of two neural networks, it is an NP-hard problem to make error ambiguity decomposition in the adversarial ensemble learning. The main challenge is that the adversarial example is heavily dependent on the intrinsic structure of neural network, and diversity analysis is relevant to intrinsic structure, model prediction and adversarial perturbation simultaneously. This is different from traditional diversity analysis on model predictions (Zhou, 2012; Wood et al., 2023). We also notice the importance of structure diversity for decision trees (Sun & Zhou, 2018), whereas it remains open for neural networks.

The proof idea involves the reduction of 3-SAT problem. Specially, we consider a 3-SAT problem with k clauses, and each clause is a disjunction of 3 literals. We construct a neural network $g(\mathbf{x})$ with $\Theta(k)$ layers and $\Theta(k)$ width, where each literal can be regarded as an input node of $g(\mathbf{x})$, and each disjunction and conjunction can be replaced with the max and min operators, respectively. The max and min operators can be constructed by ReLU activation. We construct two neural networks $f_1(\mathbf{x}, x_0) = \min(g(\mathbf{x}), x_0)$ and $f_2(\mathbf{x}, x_0) = \min(g(\mathbf{x}), -x_0)$ by adding an auxiliary variable x_0 . The detailed proof is given in Appendix A, which is partially motivated from previous work on l_1 or l_∞ norm (Katz et al., 2017; Weng et al., 2018), while our work generalizes to l_p norm with $p = 1, 2, \dots, \infty$.

3.2. Diversity decomposition w.r.t. squared loss

Previous ensemble methods generally take the first-order approximation for adversarial loss function, and focus on gradient diversity (Kariyappa & Qureshi, 2019; Dabouei et al., 2020; Huang et al., 2021). This is partially because of the NP-hardness for adversarial diversity in Theorem 3.1. Following the first-order Taylor approximation, we have

$$\bar{f}(\mathbf{x}+\delta) = \sum_{j=1}^m \frac{f_j(\mathbf{x}+\delta)}{m} \approx \sum_{j=1}^m \frac{f_j(\mathbf{x})}{m} + \sum_{j=1}^m \frac{\nabla f_j(\mathbf{x})^T \delta}{m}.$$

We now present the first decomposition for the adversarial ensemble loss w.r.t. squared loss as follows:

Theorem 3.2. For ensemble $\bar{f} = \sum_{j=1}^m f_j/m$, we present the decomposition of adversarial ensemble loss under the first-order approximation over example (\mathbf{x}, y) as follows:

$$\begin{aligned} \max_{\delta \in \Delta_p^\epsilon} \{(\bar{f}(\mathbf{x}+\delta) - y)^2\} &= \underbrace{\sum_{j=1}^m \frac{\max_{\delta \in \Delta_p^\epsilon} \{(f_j(\mathbf{x}+\delta) - y)^2\}}{m}}_{\text{average of individual adversarial losses}} \\ &- \underbrace{\sum_{j=1}^m \frac{(f_j(\mathbf{x}) - \bar{f}(\mathbf{x}))^2}{m}}_{\text{prediction diversity}} - \underbrace{\epsilon^2 \sum_{j=1}^m \frac{\|\nabla f_j(\mathbf{x})\|_q^2 - \|\nabla \bar{f}(\mathbf{x})\|_q^2}{m}}_{\text{gradient diversity}} \end{aligned}$$

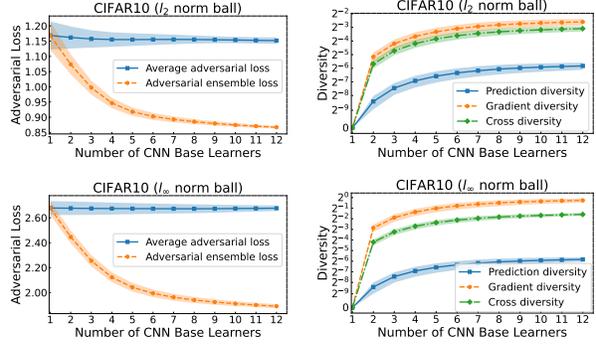


Figure 1. An illustration of diversity decomposition in Theorem 3.2 on dataset CIFAR10 of perturbation balls with l_2 and l_∞ norm, respectively. Here, we consider CNNs as base learners.

$$- 2\epsilon \underbrace{\sum_{j=1}^m \frac{\|\nabla f_j(\mathbf{x})\|_q |f_j(\mathbf{x}) - y| - \|\nabla \bar{f}(\mathbf{x})\|_q |\bar{f}(\mathbf{x}) - y|}{m}}_{\text{cross diversity}}$$

where $1/p + 1/q = 1$.

In this theorem, the adversarial ensemble loss is decomposed into the average of individual adversarial losses, *prediction diversity*, *gradient diversity* and *cross diversity*. Gradient diversity measures the dispersion of gradients concerning the mean of gradients w.r.t. l_q norm, and it exactly becomes the variance as for $q = 2$. Prediction diversity can be understood as the traditional diversity such as ambiguity (Krogh & Vedelsby, 1994; Zhou, 2012). Cross diversity can be viewed as a cross between functional outputs and gradients. Gradient and cross diversities are highly relevant to intrinsic structures of functions and functional outputs. The detailed proof is given in Appendix B.1.

From Theorem 3.2, it is also observable that gradient and cross diversities take more important roles on adversarial ensemble learning as for larger ϵ (i.e., radius of perturbation ball), yet prediction diversity dominates as for smaller ϵ . Also, the decomposition is relevant to the l_p -norm distance of perturbation ball. Thus, the characterization of diversities in adversarial ensemble learning is more complicated than that of traditional non-adversarial ensemble learning.

Figure 1 presents an intuitive illustration for the diversity decomposition of Theorem 3.2. Here, we consider dataset CIFAR10 with two classes, and focus on the perturbation ball with the popular l_2 and l_∞ norm. More experiment details are given in Appendix B.2, and we try to understand the trends of loss functions and diversities as the number of CNN-base learners increases.

As can be seen from Figure 1, we have smaller adversarial ensemble loss as for larger prediction, gradient and cross diversities w.r.t. l_2 and l_∞ norm perturbation balls, when we

keep average of individual adversarial losses stable. Generally, gradient and cross diversities are larger than prediction diversity and take more important roles. This is nicely in accordance with our Theorem 3.2, and diversity empirically plays an important role in adversarial ensemble learning.

We focus on the first-order approximation as in previous adversarial ensemble methods (Kariyappa & Qureshi, 2019; Dabouei et al., 2020; Huang et al., 2021), and it is interesting to explore the second-order (or higher-order) approximation from gradients and Hessian matrices. The main challenge is how to obtain the closed-form solution of adversarial loss via second-order approximation, because it is relevant to the roots of high-order polynomials (Forsythe & Golub, 1965; Moré & Sorensen, 1983; Fortin & Wolkowicz, 2004). More discussions on this issue are given in Appendix B.3.

Relevant to previous adversarial ensemble methods

Most previous ensemble methods consider the first-order approximation of adversarial loss functions (Kariyappa & Qureshi, 2019; Dabouei et al., 2020; Yang et al., 2021; Bogun et al., 2022), and the diversity is measured by the average of cos values over pairs of gradients of base learner. Specifically, the diversity w.r.t. instance \mathbf{x} is given by

$$\frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j=i+1}^m \cos(\nabla f_i(\mathbf{x}), \nabla f_j(\mathbf{x})),$$

where $\cos(\nabla f_i(\mathbf{x}), \nabla f_j(\mathbf{x}))$ denotes the cos value of the angle between $\nabla f_i(\mathbf{x})$ and $\nabla f_j(\mathbf{x})$. We could derive the relationship between our gradient diversity and previous diversity via the cos functions as follows:

$$\begin{aligned} \text{Gradient diversity} &= \frac{m\epsilon^2 - \epsilon^2}{m^2} \sum_{i=1}^m \|\nabla f_i(\mathbf{x})\|_2^2 \\ &- \frac{\epsilon^2}{m^2} \sum_{i \neq j} \|\nabla f_i(\mathbf{x})\|_2 \|\nabla f_j(\mathbf{x})\|_2 \cos(\nabla f_i(\mathbf{x}), \nabla f_j(\mathbf{x})). \end{aligned}$$

It is feasible to enlarge gradient diversity by decreasing cos functions, which is nicely in accordance with previous work (Dabouei et al., 2020; Bogun et al., 2022). Meanwhile, it is noteworthy of other important factors on gradient diversity such as gradient norm, rather than only one factor, which can be shown by following examples.

Example 1. There exist two ensembles of the same averages of cos values and individual adversarial losses, but with different adversarial ensemble losses.

Proof. We focus on 2-dimensional instance space $\mathcal{X} \subseteq \mathbb{R}^2$ and label space $\mathcal{Y} \subseteq \mathbb{R}$, and consider

$$\begin{aligned} f_1(\mathbf{x}) &= x_1 + x_2, & f_2(\mathbf{x}) &= x_1 - 3x_2, \\ f_3(\mathbf{x}) &= abx_1 + bx_2, & f_4(\mathbf{x}) &= abx_1 - bx_2, \end{aligned}$$

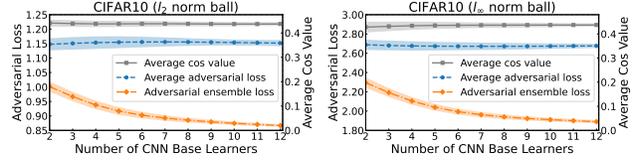


Figure 2. An illustration of the influence of average cos value on dataset CIFAR10 of perturbation balls with l_2 and l_∞ norm, respectively. Here, we consider CNNs as base learners.

where $a = (\sqrt{5}-1)/2$ and $b = (\sqrt{6}+1)/(a + \sqrt{5}a)$. We study two ensembles: one ensemble of f_1 and f_2 ; the other ensemble of f_3 and f_4 . For example $(\mathbf{x}, y) = ([1, 0], 1)$ and perturbation set $\Delta = \{\delta : \|\delta\|_2 \leq 2\}$, we have

$$\cos(\nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x})) = \cos(\nabla f_3(\mathbf{x}), \nabla f_4(\mathbf{x})),$$

also with the same average of individual adversarial losses

$$\sum_{i=1}^2 \max_{\delta \in \Delta} \frac{(f_i(\mathbf{x} + \delta) - y)^2}{2} = \sum_{i=3}^4 \max_{\delta \in \Delta} \frac{(f_i(\mathbf{x} + \delta) - y)^2}{2}.$$

However, the adversarial ensemble losses are different from

$$\begin{aligned} \max_{\delta \in \Delta} (f_1(\mathbf{x} + \delta)/2 + f_2(\mathbf{x} + \delta)/2 - y)^2 &= 8, \\ \max_{\delta \in \Delta} (f_3(\mathbf{x} + \delta)/2 + f_4(\mathbf{x} + \delta)/2 - y)^2 &\approx 7.2. \end{aligned}$$

Here, we consider the l_2 norm in perturbation set Δ , and similar analysis could be made for l_p -norm. More details are presented in Appendix B.4. \square

In addition to Example 1, we could also present empirical studies on adversarial ensemble losses versus the cos values over dataset CIFAR10, and the experiment details are given in Appendix B.2. Figure 2 shows the curves of adversarial ensemble loss, the averages of cos values and individual adversarial losses with l_2 and l_∞ norm.

From Figure 2, it is clear that adversarial ensemble losses keep decreasing when we increase number of CNN-base learners, whereas it almost remains constant for the averages of cos values and individual adversarial losses. Therefore, it is not sufficient to characterize the diversity of adversarial ensemble by merely considering the average of cos values as in (Dabouei et al., 2020; Bogun et al., 2022).

3.3. Diversity decomposition w.r.t. cross-entropy loss

We study the decomposition of adversarial ensemble w.r.t. cross-entropy loss for binary classification, and also follow the first-order approximation $f(\mathbf{x} + \delta) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta$. For base learners with logit outputs (i.e., the logarithm of the ratio of probabilities), we have the probability of the positive class over $\mathbf{x} + \delta$ as follows:

$$p_f(\mathbf{x} + \delta) = (1 + \exp(-(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta)))^{-1}. \quad (1)$$

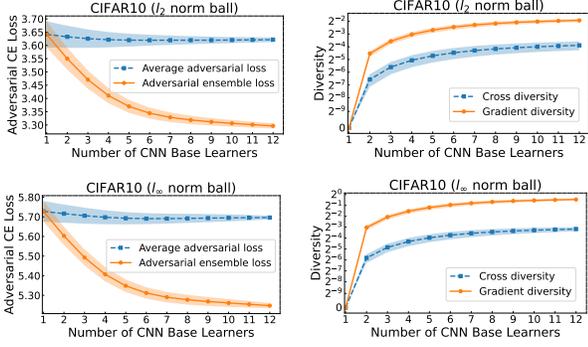


Figure 3. An illustration of diversity decomposition in Theorem 3.4 on dataset CIFAR10 of perturbation balls with l_2 and l_∞ norm, respectively. Here, we consider CNNs as base learners.

The cross-entropy loss over $(\mathbf{x} + \boldsymbol{\delta}, y)$ is given by

$$\begin{aligned} \ell(f(\mathbf{x} + \boldsymbol{\delta}), y) &= -y(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \boldsymbol{\delta}) \\ &\quad + \ln(1 + \exp(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \boldsymbol{\delta})), \end{aligned}$$

as in (Bishop & Nasrabadi, 2006), and recall that

$$\boldsymbol{\delta}_f^* \in \arg \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \ell(f(\mathbf{x} + \boldsymbol{\delta}), y). \quad (2)$$

We have the following closed-form solution for $p_f(\mathbf{x} + \boldsymbol{\delta}^*)$, and the detailed proof is presented in Appendix C.1.

Lemma 3.3. For function f and example (\mathbf{x}, y) , we have

$$p_f(\mathbf{x} + \boldsymbol{\delta}_f^*) = (1 + \exp(-(f(\mathbf{x}) - (2y - 1)\|\nabla f(\mathbf{x})\|_q \epsilon)))^{-1}$$

where probability of positive class $p_f(\cdot)$ and perturbation $\boldsymbol{\delta}_f^*$ are defined by Eqns. (1) and (2), respectively.

Based on this lemma, we have

Theorem 3.4. For ensemble $\bar{f} = \sum_{j=1}^m f_j/m$, we have the decomposition of adversarial ensemble loss under the first-order approximation over example (\mathbf{x}, y) as follows:

$$\begin{aligned} \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \ell(\bar{f}(\mathbf{x} + \boldsymbol{\delta}), y) &= \underbrace{\sum_{j=1}^m \frac{\max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \ell(f_j(\mathbf{x} + \boldsymbol{\delta}), y)}{m}}_{\text{average of individual adversarial losses}} \\ &\quad - r \underbrace{\sum_{j=1}^m \frac{\|\nabla f_j(\mathbf{x})\|_q - \|\nabla \bar{f}(\mathbf{x})\|_q}{m}}_{\text{gradient diversity}} - \underbrace{\sum_{j=1}^m \frac{KL(p_{\bar{f}}(\tilde{\mathbf{x}}_{\bar{f}}), p_{f_j}(\tilde{\mathbf{x}}_{f_j}))}{m}}_{\text{cross diversity}} \end{aligned}$$

where $1/p + 1/q = 1$, $r = \epsilon(y - p_{\bar{f}}(\tilde{\mathbf{x}}_{\bar{f}}))(2y - 1)$, $p_f(\cdot)$ is defined by Eqn. (1), and $\tilde{\mathbf{x}}_f = \mathbf{x} + \boldsymbol{\delta}_f^*$ is the adversarial example with $\boldsymbol{\delta}_f^*$ from Eqn. (2).

In this theorem, the adversarial ensemble loss is decomposed into the average of individual adversarial losses, *gradient*

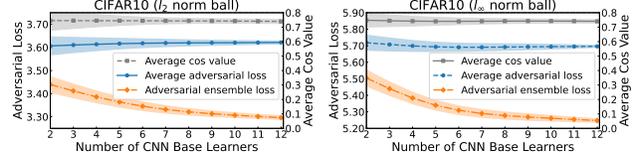


Figure 4. An illustration of influence of the average cos value on dataset CIFAR10 of perturbation balls w.r.t. l_2 and l_∞ norm, respectively. Here, we consider CNNs as base learners.

diversity and *cross diversity*. The cross diversity shows the diversity of base learners according to their KL divergence to ensemble \bar{f} , which is relevant to function outputs and gradients, as shown in Lemma 3.3. It is evident that the decomposition of Theorem 3.4 is quite different from that of Theorem 3.2 because of different loss functions. The detailed proof is presented in Appendix C.2.

We also present Figure 3 to illustrate the decomposition of Theorem 3.4, and some experimental details are given in Appendix B.2. As can be seen, adversarial ensemble loss gets smaller as for larger cross and gradient diversities w.r.t. l_2 and l_∞ norm perturbation balls, when we maintain the average of individual adversarial losses stable. In addition, it is not sufficient to characterize the diversity of adversarial ensemble by only merely considering the average of cos values as done in (Dabouei et al., 2020; Bogun et al., 2022). We also present some empirical studies on the influence of average cos values on adversarial ensemble loss shown in Figure 4, and more details are given in Appendix C.3.

4. Our AdvE_{OAP} Method

Motivated from our theoretical analysis in Theorem 3.4, we develop a robust ensemble method for multi-class learning, and the basic idea is to train multiple deep neural networks for adversarial ensemble via a regularization by considering gradient diversity and cross diversity simultaneously.

For multi-class learning with κ classes, the base learner $\mathbf{f} = (f_1, f_2, \dots, f_\kappa): \mathcal{X} \rightarrow \mathbb{R}^\kappa$ maps each instance to a κ -dimensional logit vector. The predicted probability vector $\mathbf{p}_f(\mathbf{x}) = (p_{f_1}(\mathbf{x}), p_{f_2}(\mathbf{x}), \dots, p_{f_\kappa}(\mathbf{x}))$ can be calculated from $\mathbf{f}(\mathbf{x})$ via softmax function as follows

$$p_{f_k}(\mathbf{x}) = \frac{\exp(f_k(\mathbf{x}))}{\sum_{j=1}^{\kappa} \exp(f_j(\mathbf{x}))}. \quad (3)$$

This section also focuses on the first-order approximation $\mathbf{f}(\mathbf{x} + \boldsymbol{\delta}) \approx \mathbf{f}(\mathbf{x}) + \mathbf{J}_f(\mathbf{x})\boldsymbol{\delta}$, where $\mathbf{J}_f(\mathbf{x})$ is the Jacobi matrix of \mathbf{f} (Rudin et al., 1964). Denote by

$$\boldsymbol{\delta}_f^* \in \arg \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \ell(\mathbf{f}(\mathbf{x} + \boldsymbol{\delta}), y). \quad (4)$$

For m base learners $\mathbf{f}_1, \dots, \mathbf{f}_m$ with $\mathbf{f}_j = (f_{j,1}, \dots, f_{j,\kappa})$,

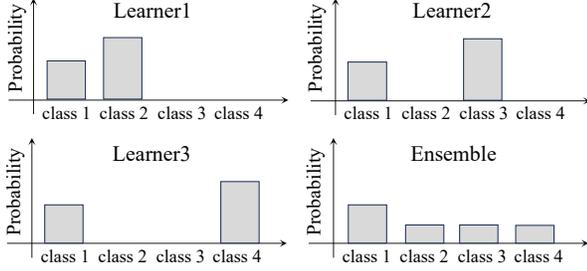


Figure 5. An illustration of orthogonality by optimizing Eqn. (5).

we have cross diversity for multi-classification

$$\text{Cross diversity} = \sum_{j=1}^m \frac{\text{KL}(\mathbf{p}_{\bar{f}}(\mathbf{x} + \delta_{\bar{f}}^*), \mathbf{p}_{f_j}(\mathbf{x} + \delta_{f_j}^*))}{m},$$

where $\mathbf{p}_f(\cdot)$ and δ_f^* are given by Eqns. (3)-(4), respectively. From some algebraic derivations in Appendix D.1, we have

$$\begin{aligned} \text{Gradient diversity} &= \sum_{j=1}^m \frac{\langle \mathbf{r}, \mathbf{J}_{f_j}(\mathbf{x})\delta_{f_j}^* - \mathbf{J}_{\bar{f}}(\mathbf{x})\delta_{\bar{f}}^* \rangle}{m} \\ &= \sum_{j=1}^m \frac{\langle \mathbf{r}, f_j(\mathbf{x} + \delta_{f_j}^*) - \bar{f}(\mathbf{x} + \delta_{\bar{f}}^*) \rangle}{m} \end{aligned}$$

where $\mathbf{r} = \mathbf{p}_{\bar{f}}(\mathbf{x} + \delta_{\bar{f}}^*) - \mathbf{e}_y$.

For multi-class learning, we could simultaneously improve cross diversity and gradient diversity by diversifying the output predictions of adversarial examples of base learners. This motivates us to develop a new ensemble algorithm via orthogonal adversarial predictions, i.e., we orthogonalize the outputs over adversarial examples for base learners to improve diversity. The orthogonal idea is inspired by (Pang et al., 2019), which is limited only on clean examples, while our work generalizes to adversarial examples.

Specifically, we introduce a regularization for orthogonal adversarial output predictions of f_1, \dots, f_m as

$$\Gamma_\alpha(\mathbf{x}, y) = H \left(\sum_{j=1}^m \frac{\tilde{\mathbf{p}}_{f_j}(\mathbf{x} + \delta_{f_j}^*)}{m} \right) + \alpha \log(\det(\mathbf{A}^T \mathbf{A})),$$

where $H(\cdot)$ is the function of information entropy and

$$\mathbf{A} = [\tilde{\mathbf{p}}_{f_1}(\mathbf{x} + \delta_{f_1}^*), \dots, \tilde{\mathbf{p}}_{f_m}(\mathbf{x} + \delta_{f_m}^*)] \in \mathbb{R}^{(K-1) \times m}.$$

Here, $\tilde{\mathbf{p}}_f(\cdot)$ is an $(K-1)$ -dimensional vector obtained by removing the y -th element of $\mathbf{p}_f(\cdot)$ and normalizing with l_1 norm, and $\det(\mathbf{A}^T \mathbf{A})$ shows the square of volume of polytope spanned by $\tilde{\mathbf{p}}_{f_1}(\mathbf{x} + \delta_{f_1}^*), \dots, \tilde{\mathbf{p}}_{f_m}(\mathbf{x} + \delta_{f_m}^*)$ as in (Bernstein, 2009). Intuitively, the orthogonal probability vectors $\tilde{\mathbf{p}}_{f_1}(\mathbf{x} + \delta_{f_1}^*), \dots, \tilde{\mathbf{p}}_{f_m}(\mathbf{x} + \delta_{f_m}^*)$ could yield

Algorithm 1 The AdvE_{OAP} method

Input: training dataset S , number of base learners m , learning rate η and the SGD iterations T

Initialize: base learners f_1, \dots, f_m

for $t = 1$ **to** T **do**

Partition training dataset S into batches S_1, \dots, S_b

for $k \in [b]$ **do**

Generate adversarial examples for every base learner and example by the PGD-attack

for $j = 1$ **to** m **do**

Update learner f_j by gradient descent in Eqn. (5)

end for

end for

end for

Output: Ensemble model $\bar{f} = \sum_{j=1}^m f_j/m$

larger volume and entropy in $\Gamma_\alpha(\mathbf{x}, y)$, and hence improve gradient diversity and cross diversity simultaneously.

Given training sample $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, we present the final objective optimization as

$$\sum_{i=1}^n \sum_{j=1}^m \max_{\delta \in \Delta_p^c} \frac{\ell(f_j(\mathbf{x}_i + \delta), y_i) - \lambda \Gamma_\alpha(\mathbf{x}_i, y_i)}{mn}, \quad (5)$$

where λ is a hyper-parameter to tradeoff adversarial loss and regularization $\Gamma_\alpha(\cdot)$. We could obtain mutual orthogonal base learners f_1, \dots, f_m from the optimization of Eqn. (5), and the details are given in Appendix D.2.

Figure 5 gives an illustration for the orthogonality of base learners. Here are three base learners in multi-class learning of 4 classes, and the ground-truth class is class 1. The output predictions are mutually orthogonal except for the ground-truth class 1. Therefore, the ensemble of three base learners could make correct prediction robustly even if three base learners are misled to different classes.

On the optimization of Eqn. (5), we take adversarial training method with stochastic gradient descent from (Madry et al., 2018). We generate adversarial examples w.r.t. all base learners by PGD-attack, and calculate gradients to update the parameters for each base learner of deep neural network. Algorithm 1 presents the detailed description for our Adversarial Ensemble training with Orthogonal Adversarial Predictions, which is short for AdvE_{OAP}. For Algorithm 1, the time complexity takes m -times as that of training a single neural network adversarially. In addition, it takes $O(m^3)$ computational cost to calculate the regularization and gradients. More details are given in Appendix D.3

Table 1. Comparison of classification accuracies (mean \pm std %) over adversarial examples generated by different adversarial attacks. \bullet / \circ indicates that our AdvE_{OAP} is significantly better/worse than the corresponding methods (pair-wise t -test at 95% significance level).

	Methods	FGSM	PGD10	PGD20	PGD40	AutoPGD	MORA	AutoAttack
MNIST	Our AdvE _{OAP}	96.413 \pm 0.124	95.676 \pm 0.039	95.648 \pm 0.056	95.504 \pm 0.098	94.907 \pm 0.187	94.672 \pm 0.209	94.072 \pm 0.423
	GAL	10.271 \pm 0.709 \bullet	00.001 \pm 0.002 \bullet	00.002 \pm 0.005 \bullet	00.000 \pm 0.000 \bullet	00.000 \pm 0.000 \bullet	00.004 \pm 0.005 \bullet	00.000 \pm 0.000 \bullet
	ADP	10.888 \pm 1.931 \bullet	00.000 \pm 0.000 \bullet	00.004 \pm 0.006 \bullet	00.000 \pm 0.000 \bullet			
	AdvADP	95.333 \pm 0.098 \bullet	95.119 \pm 0.061	95.059 \pm 0.089	93.032 \pm 0.191 \bullet	93.298 \pm 0.109 \bullet	93.172 \pm 0.063 \bullet	92.948 \pm 0.062 \bullet
	DVERGE	74.903 \pm 1.059 \bullet	37.506 \pm 4.292 \bullet	34.807 \pm 4.981 \bullet	05.846 \pm 2.573 \bullet	00.001 \pm 0.002 \bullet	00.299 \pm 0.361 \bullet	00.000 \pm 0.000 \bullet
	PDD	10.446 \pm 1.293 \bullet	06.041 \pm 2.640 \bullet	04.059 \pm 2.904 \bullet	02.256 \pm 2.792 \bullet	01.163 \pm 1.295 \bullet	04.100 \pm 4.006 \bullet	00.000 \pm 0.000 \bullet
	TRS	91.044 \pm 0.892 \bullet	86.544 \pm 1.300 \bullet	86.709 \pm 1.014 \bullet	85.954 \pm 3.283 \bullet	80.902 \pm 4.905 \bullet	79.628 \pm 5.471 \bullet	78.615 \pm 5.997 \bullet
	iGAT _{ADP}	83.814 \pm 2.408 \bullet	79.892 \pm 1.410 \bullet	79.281 \pm 1.732 \bullet	79.778 \pm 2.028 \bullet	59.357 \pm 4.559 \bullet	50.208 \pm 6.439 \bullet	48.589 \pm 5.178 \bullet
F-MNIST	Our AdvE _{OAP}	82.743 \pm 0.263	81.770 \pm 0.245	81.752 \pm 0.248	80.543 \pm 0.290	81.467 \pm 0.079	80.643 \pm 0.332	80.506 \pm 0.327
	GAL	15.404 \pm 6.709 \bullet	00.352 \pm 0.447 \bullet	00.170 \pm 0.215 \bullet	00.315 \pm 0.306 \bullet	00.000 \pm 0.000 \bullet	00.009 \pm 0.005 \bullet	00.000 \pm 0.000 \bullet
	ADP	22.287 \pm 1.741 \bullet	00.001 \pm 0.002 \bullet	00.001 \pm 0.002 \bullet	00.000 \pm 0.000 \bullet	00.000 \pm 0.000 \bullet	00.009 \pm 0.005 \bullet	00.000 \pm 0.000 \bullet
	AdvADP	82.728 \pm 0.331	79.167 \pm 0.338 \bullet	79.074 \pm 0.387 \bullet	80.080 \pm 0.408 \bullet	78.461 \pm 0.372 \bullet	77.833 \pm 0.313 \bullet	77.732 \pm 0.245 \bullet
	DVERGE	48.242 \pm 1.536 \bullet	27.140 \pm 2.878 \bullet	25.846 \pm 3.019 \bullet	32.652 \pm 3.033 \bullet	17.222 \pm 4.786 \bullet	20.126 \pm 4.605 \bullet	15.035 \pm 5.560 \bullet
	PDD	29.936 \pm 4.991 \bullet	18.740 \pm 4.925 \bullet	17.958 \pm 4.873 \bullet	19.043 \pm 4.524 \bullet	12.161 \pm 5.521 \bullet	14.408 \pm 3.617 \bullet	00.319 \pm 0.466 \bullet
	TRS	70.767 \pm 0.466 \bullet	69.330 \pm 0.125 \bullet	69.285 \pm 0.105 \bullet	67.641 \pm 0.440 \bullet	68.719 \pm 0.098 \bullet	67.800 \pm 0.183 \bullet	67.250 \pm 0.304 \bullet
	iGAT _{ADP}	66.278 \pm 2.051 \bullet	63.139 \pm 0.236 \bullet	62.967 \pm 0.149 \bullet	62.882 \pm 0.379 \bullet	61.920 \pm 0.083 \bullet	51.869 \pm 4.773 \bullet	48.750 \pm 5.763 \bullet
CIFAR10	Our AdvE _{OAP}	55.718 \pm 0.245	53.076 \pm 0.249	52.996 \pm 0.255	52.903 \pm 0.295	51.997 \pm 0.234	48.318 \pm 0.065	47.884 \pm 0.060
	GAL	12.370 \pm 2.959 \bullet	00.007 \pm 0.013 \bullet	00.000 \pm 0.000 \bullet	00.000 \pm 0.000 \bullet	00.000 \pm 0.000 \bullet	00.002 \pm 0.004 \bullet	00.000 \pm 0.000 \bullet
	ADP	23.100 \pm 0.757 \bullet	00.008 \pm 0.010 \bullet	00.001 \pm 0.003 \bullet	00.000 \pm 0.000 \bullet	00.000 \pm 0.000 \bullet	00.001 \pm 0.003 \bullet	00.000 \pm 0.000 \bullet
	AdvADP	55.478 \pm 0.214	47.116 \pm 0.278 \bullet	46.802 \pm 0.284 \bullet	46.573 \pm 0.332 \bullet	44.192 \pm 0.216 \bullet	42.904 \pm 0.235 \bullet	42.174 \pm 0.198 \bullet
	DVERGE	28.536 \pm 0.882 \bullet	05.358 \pm 0.344 \bullet	04.830 \pm 0.343 \bullet	04.690 \pm 0.269 \bullet	02.246 \pm 0.155 \bullet	02.868 \pm 0.178 \bullet	01.748 \pm 0.164 \bullet
	PDD	23.750 \pm 4.005 \bullet	14.116 \pm 4.596 \bullet	14.896 \pm 4.596 \bullet	17.400 \pm 2.710 \bullet	05.570 \pm 2.365 \bullet	09.526 \pm 4.345 \bullet	01.000 \pm 0.626 \bullet
	TRS	39.350 \pm 0.402 \bullet	37.548 \pm 0.431 \bullet	37.453 \pm 0.440 \bullet	37.263 \pm 0.259 \bullet	36.690 \pm 0.440 \bullet	32.828 \pm 0.433 \bullet	32.588 \pm 0.433 \bullet
	iGAT _{ADP}	19.990 \pm 0.509 \bullet	18.075 \pm 0.035 \bullet	18.000 \pm 0.056 \bullet	13.707 \pm 2.596 \bullet	17.260 \pm 0.099 \bullet	12.365 \pm 1.054 \bullet	12.090 \pm 1.103 \bullet

5. Experiments

We conduct experiments on three datasets²: MNIST of 70000 images and 784 dimensions, F-MNIST of 70000 images and 784 dimensions, and CIFAR10 of 60000 images and 3072 dimensions. Three datasets have been well-studied in previous works (Strauss et al., 2017; Kariyappa & Qureshi, 2019; Yang et al., 2021; Deng & Mu, 2024). We compare our method with the state-of-the-art methods on adversarial ensemble learning as follows:

- GAL: Non-adversarial training via the diversity of cos values of gradients (Kariyappa & Qureshi, 2019);
- ADP: Non-adversarial training via the diversity of the orthogonality of predictions (Pang et al., 2019);
- AdvADP: Adversarial training via the diversity of the orthogonality of predictions (Pang et al., 2019);
- DVERGE: Adversarial training on the exchange of adversarial examples in base learners to diversify the adversarial vulnerability (Yang et al., 2020);

- PDD: Non-adversarial training to diversify the feature representations via dropouts (Huang et al., 2021);
- TRS: GAL by preserving the smoothness of base learners (Yang et al., 2021);
- iGAT_{ADP}: AdvADP by allocating globally adversarial examples to base learners (Deng & Mu, 2024).

For all datasets, the perturbation size is set as 0.2, 0.05 and 0.03 under l_∞ -norm ball, respectively, as done in (Croce & Hein, 2020; Deng & Mu, 2024). For all methods, we select ResNet20 as base learners with learners number as 3, 3 and 8 for MNIST, F-MNIST and CIFAR10, respectively. More settings are given in Appendix E.1. All experiments are performed on a server with 64 CPU cores (2 Intel Xeon Gold 6430 CPUs) and NVIDIA GeForce RTX 4090 GPU, running Ubuntu 24.04 with 1TB main memory.

5.1. Performance under adversarial attacks

We take accuracy to measure performance on adversarial examples generated by seven popular adversarial attacks, i.e., FGSM (Goodfellow et al., 2015), PGD10, PGD20,

²Download from <https://paperswithcode.com/dataset>.

Table 2. Comparison of classification accuracies (mean \pm std %) over adversarial examples generated by EOT and BPDA attacks. \bullet / \circ indicates that our AdvE_{OAP} is significantly better/worse than the corresponding methods (pair-wise t -test at 95% significance level).

Attacks	Datasets	Our AdvE _{OAP}	GAL	AdvADP	PDD	DVERGE	TRS	iGAT _{ADP}
EOT	MNIST	88.116 \pm 0.316	01.178 \pm 1.958 \bullet	85.724 \pm 0.163 \bullet	03.913 \pm 5.471 \bullet	38.358 \pm 3.147 \bullet	75.535 \pm 4.981 \bullet	72.721 \pm 1.258 \bullet
	F-MNIST	62.416 \pm 0.778	02.201 \pm 1.886 \bullet	61.137 \pm 0.547 \bullet	08.323 \pm 1.815 \bullet	35.162 \pm 2.390 \bullet	56.592 \pm 1.110 \bullet	57.893 \pm 2.147 \bullet
	CIFAT10	42.003 \pm 0.504	01.490 \pm 0.179 \bullet	40.787 \pm 0.238 \bullet	08.535 \pm 0.006 \bullet	20.371 \pm 0.550 \bullet	31.629 \pm 1.156 \bullet	16.408 \pm 2.697 \bullet
BPDA ₁	MNIST	95.382 \pm 0.079	00.002 \pm 0.005 \bullet	92.713 \pm 0.137 \bullet	02.757 \pm 2.852 \bullet	14.187 \pm 7.868 \bullet	85.024 \pm 3.758 \bullet	71.763 \pm 3.041 \bullet
	F-MNIST	80.080 \pm 0.256	00.352 \pm 0.447 \bullet	79.074 \pm 0.387 \bullet	17.958 \pm 4.873 \bullet	25.846 \pm 3.019 \bullet	66.409 \pm 0.450 \bullet	62.967 \pm 0.149 \bullet
	CIFAT10	49.530 \pm 0.050	00.000 \pm 0.000 \bullet	44.687 \pm 0.286 \bullet	14.896 \pm 4.596 \bullet	04.730 \pm 0.343 \bullet	33.417 \pm 0.310 \bullet	10.693 \pm 1.776 \bullet
BPDA ₂	MNIST	95.676 \pm 0.104	00.000 \pm 0.000 \bullet	93.646 \pm 0.133 \bullet	12.398 \pm 4.848 \bullet	19.385 \pm 8.467 \bullet	86.783 \pm 3.416 \bullet	84.261 \pm 1.767 \bullet
	F-MNIST	80.832 \pm 0.311	01.098 \pm 0.243 \bullet	80.832 \pm 0.380	30.152 \pm 8.392 \bullet	40.672 \pm 2.473 \bullet	68.032 \pm 0.364 \bullet	73.444 \pm 0.657 \bullet
	CIFAT10	53.092 \pm 0.230	00.000 \pm 0.000 \bullet	47.858 \pm 0.174 \bullet	25.277 \pm 4.449 \bullet	05.608 \pm 0.343 \bullet	37.383 \pm 0.355 \bullet	14.500 \pm 2.608 \bullet
BPDA ₃	MNIST	96.132 \pm 0.071	00.000 \pm 0.000 \bullet	94.857 \pm 0.085 \bullet	09.165 \pm 0.941 \bullet	57.491 \pm 1.640 \bullet	88.993 \pm 2.873 \bullet	91.126 \pm 0.672 \bullet
	F-MNIST	81.309 \pm 0.409	00.343 \pm 0.206 \bullet	82.106 \pm 0.468	30.265 \pm 9.179 \bullet	52.157 \pm 1.658 \bullet	68.900 \pm 0.330 \bullet	75.948 \pm 0.619 \bullet
	CIFAT10	52.440 \pm 0.226	00.000 \pm 0.000 \bullet	46.240 \pm 0.252 \bullet	25.690 \pm 3.737 \bullet	05.347 \pm 0.382 \bullet	36.480 \pm 0.261 \bullet	13.600 \pm 2.548 \bullet
BPDA ₄	MNIST	94.596 \pm 0.135	00.022 \pm 0.020 \bullet	92.830 \pm 0.244 \bullet	04.941 \pm 4.538 \bullet	02.635 \pm 1.482 \bullet	81.595 \pm 2.255 \bullet	44.669 \pm 7.900 \bullet
	F-MNIST	80.241 \pm 0.247	00.026 \pm 0.045 \bullet	80.757 \pm 0.386	05.706 \pm 7.595 \bullet	50.324 \pm 4.693 \bullet	67.526 \pm 0.467 \bullet	68.174 \pm 0.462 \bullet
	CIFAT10	53.383 \pm 0.110	00.273 \pm 0.073 \bullet	53.397 \pm 0.160	12.650 \pm 1.642 \bullet	26.573 \pm 0.046 \bullet	40.507 \pm 0.399 \bullet	14.410 \pm 1.161 \bullet

Table 3. Comparison of accuracies (mean \pm std%) for our AdvE_{OAP} with and without regularization $\Gamma(\cdot)$ under the APGD attack.

Our AdvE _{OAP}	MNIST	F-MNIST	CIFAR10
without $\Gamma_\alpha(\cdot)$	93.513 \pm 0.046	80.317 \pm 0.159	50.833 \pm 0.102
with $\Gamma_\alpha(\cdot)$	94.907 \pm 0.187	81.202 \pm 0.311	51.997 \pm 0.234
Improvement	1.394 \pm 0.149 \uparrow	1.150 \pm 0.140 \uparrow	1.163 \pm 0.257 \uparrow

PGD40 (Madry et al., 2018), AutoPGD (Croce & Hein, 2020), MORA (Gao et al., 2022) and AutoAttack (Croce & Hein, 2020). All methods are evaluated over 50 runs with different random initializations, as summarized in Table 1.

From Table 1, it is clear that our AdvE_{OAP} method achieves significantly better performance than GAL, PDD and TRS, since it wins at most times and never loses. This is because such methods merely focus on \cos values as the diversity measure, yet ignore other factors such as cross diversity and individual adversarial losses, which is consistent with the ensemble decomposition as in Theorem 3.4.

Our AdvE_{OAP} is also better than DVERGE and iGAT_{ADP}, since DVERGE exchanges adversarial examples of base learners and iGAT_{ADP} allocates adversarial examples of ensemble, without the consideration of diversities of base learners. Our AdvE_{OAP} outperforms ADP and AdvADP, since such methods consider diversity over clean examples, rather than adversarial examples. It is quite different to study diversity in the adversarial ensemble learning, which is heavily relevant to intrinsic structures and output predictions of models simultaneously as in Theorem 3.1.

We further study two additional adversarial attacks BPDA (Athalye et al., 2018a) and EOT (Athalye et al., 2018b),

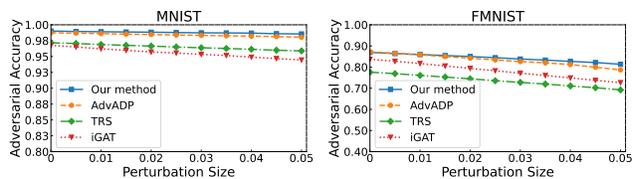


Figure 6. Influence of perturbation sizes under the PGD20 attack.

where EOT considers adversarial perturbations insensitive to transformations, and BPDA considers potential gradient risks and designs corresponded attacks. More details are given in Appendix E.2. We implement EOT and four BPDA attacks, and experimental comparisons are summarized in Table 2. It is obvious that our AdvE_{OAP} method achieves better performance than other adversarial ensembles, which shows the robustness of our method to gradient risks and adversarial perturbations insensitive to transformations.

We also present some ablation experiments to verify the effectiveness of the regularization $\Gamma_\alpha(\cdot)$ in Eqn. (5), which essentially considers gradient diversity and cross diversity via orthogonal adversarial predictions. Table 3 shows some experimental comparisons for our AdvE_{OAP} with and without regularization. It is clear that our method takes better performance with regularization, which nicely shows the importance of diversity on the design of ensemble methods.

We study the influence of different perturbation size ϵ over datasets MNIST and FMNIST, as shown in Figure 6. It is observable that our AdvE_{OAP} achieves better or comparable performance than other ensemble methods for different size of perturbations, in particular for larger ϵ . This shows the effectiveness of our methods for larger perturbation size.

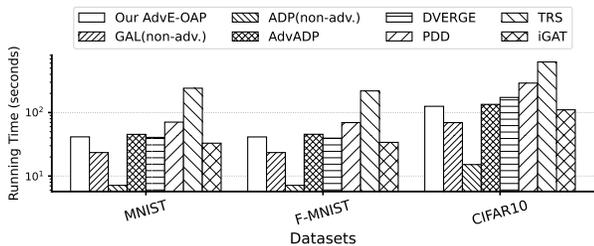


Figure 7. Comparisons of running time (seconds/epoch) for our AdvE_{OAP} and other adversarial ensemble methods.

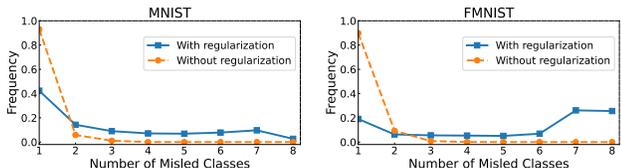


Figure 8. The frequency of number of misled classes.

We finally compare the training time of AdvE_{OAP} with other methods in Figure 7. As can be seen, our AdvE_{OAP} takes comparable training time to other adversarial-training ensembles AdvADP, DVERGE, PDD, TRS and iGAT, but with more time than the non-adversarial training methods GAL and ADP, which obviously take smaller adversarial prediction accuracy as shown in Table 1.

5.2. Orthogonality and convergence analysis

We now illustrate the orthogonality of base learners for our AdvE_{OAP}, which could mislead base learners to different classes under adversarial attacks. Figure 8 summarizes the frequency of number of misled classes with 8 base learners over two datasets MNIST and FMNIST. It is clear that base learners of our AdvE_{OAP} predict with more different classes than that of AdvE_{OAP} without regularization.

We also present the convergence analysis on adversarial ensemble loss, average of adversarial losses, as well as gradient diversity and cross diversity during the training process. Figure 9 presents the convergence curves over three datasets MNIST, FMNIST and CIFAR10. It is clear that our AdvE_{OAP} method could decrease adversarial ensemble loss, and simultaneously increase gradient diversity and cross diversity. This is nicely in accordance with our diversity decomposition in Theorem 3.4.

6. Conclusion

Diversity has always been one of the most crucial factors on the designs of ensemble methods. This work focuses on the fundamental problems of diversity in the adversarial

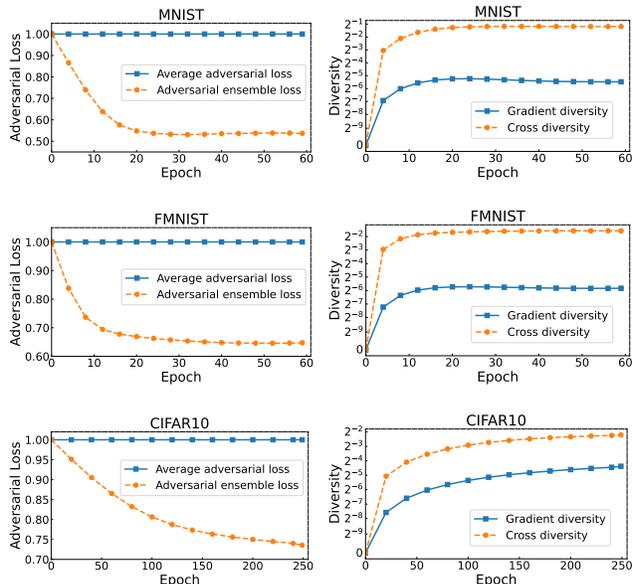


Figure 9. The curve of adversarial losses and diversities during the training process, where we normalize according to the average of individual adversarial losses.

ensemble learning. We prove the NP-Hard problem on the precise calculation of diversity for networks in adversarial ensemble learning, and give the first diversity decomposition under first-order approximation. Specifically, adversarial ensemble loss can be decomposed into average of individual adversarial losses, prediction diversity, gradient diversity and cross diversity. We consider similar diversity decomposition for classification with cross-entropy loss. Based on theoretical analysis, we develop a new ensemble method via orthogonal adversarial predictions to improve gradient and cross diversity simultaneously. An interesting future work is to explore other adversarial ensemble algorithms with better robustness and generalization from our theoretical analysis.

Acknowledgements

The authors want to thank the reviewers for their helpful comments and suggestions. This research was supported by National Key R&D Program of China (2021ZD0112802) and NSFC (62376119).

Impact Statement

This work shows the NP-Hardness of the calculation for diversity in adversarial ensemble learning and presents the first diversity decomposition with first-order approximation. It further develops a new ensemble method for adversarial defense and validates it through a series of experiments. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

References

- Abernethy, J., Awasthi, P., and Kale, S. A multiclass boosting framework for achieving fast and provable adversarial robustness. *CoRR*, abs/2103.01276, 2021.
- Andriushchenko, M., Croce, F., Flammarion, N., and Hein, M. Square attack: A query-efficient black-box adversarial attack via random search. In *Proceedings of the 16th European Conference on Computer Vision*, pp. 484–501, Virtual Event, 2020.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 274–283, Stockholm, Sweden, 2018a.
- Athalye, A., Engstrom, L., Ilyas, A., and Kwok, K. Synthesizing robust adversarial examples. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 284–293, Stockholm, Sweden, 2018b.
- Bartoldson, B. R., Diffenderfer, J., Parasyris, K., and Kailkhura, B. Adversarial robustness limits via scaling-law and human-alignment studies. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 3046–3072, Vienna, Austria, 2024.
- Bernstein, D. S. *Matrix Mathematics: Theory, Facts, and Formulas*. Princeton University Press, 2009.
- Bishop, C. M. and Nasrabadi, N. M. *Pattern Recognition and Machine Learning*, volume 4. Springer, 2006.
- Bogun, A., Kostadinov, D., and Borth, D. Saliency diversified deep ensemble for robustness to adversaries. In *Proceedings of the 36th AAAI Workshop on Adversarial Machine Learning and Beyond*, Virtual Event, 2022.
- Croce, F. and Hein, M. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 2206–2216, Vienna, Austria, 2020.
- Dabouei, A., Soleymani, S., Taherkhani, F., Dawson, J., and Nasrabadi, N. M. Exploiting joint robustness to adversarial perturbations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1122–1131, Seattle, WA, 2020.
- Deng, Y. and Mu, T. Understanding and improving ensemble adversarial defense. In *Advances in Neural Information Processing Systems 37*, pp. 58075–58087, New Orleans, LA, 2024.
- Deng, Y., Zhang, T., Lou, G., Zheng, X., Jin, J., and Han, Q.-L. Deep learning-based autonomous driving systems: A survey of attacks and defenses. *IEEE Transactions on Industrial Informatics*, 17(12):7897–7912, 2021.
- Dietterich, T. G. Ensemble methods in machine learning. In *Proceedings of the 1st International Workshop on Multiple Classifier Systems*, pp. 1–15, Cagliari, Italy, 2000.
- Finlayson, S. G., Bowers, J. D., Ito, J., Zittrain, J. L., Beam, A. L., and Kohane, I. S. Adversarial attacks on medical machine learning. *Science*, 363(6433):1287–1289, 2019.
- Forsythe, G. E. and Golub, G. H. On the stationary values of a second-degree polynomial on the unit sphere. *Journal of the Society for Industrial and Applied Mathematics*, 13(4):1050–1068, 1965.
- Fortin, C. and Wolkowicz, H. The trust region subproblem and semidefinite programming. *Optimization Methods and Software*, 19(1):41–67, 2004.
- Fursov, I., Morozov, M., Kaploukhaya, N., Kovtun, E., Rivera-Castro, R., Gusev, G., Babaev, D., Kireev, I., Zaytsev, A., and Burnaev, E. Adversarial attacks on deep models for financial transaction records. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 2868–2878, Virtual Event, 2021.
- Gao, X., Xu, C.-Z., et al. Mora: Improving ensemble robustness evaluation with model reweighing attack. In *Advances in Neural Information Processing Systems 35*, pp. 26955–26965, New Orleans, LA, 2022.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. In *Proceedings of the 3rd International Conference on Learning Representations*, San Diego, CA, 2015.
- Gowal, S., Rebuffi, S.-A., Wiles, O., Stimberg, F., Calian, D. A., and Mann, T. A. Improving robustness using generated data. In *Advances in Neural Information Processing Systems 34*, pp. 4218–4233, Virtual Event, 2021.
- Guo, J.-Q., Teng, M.-Z., Gao, W., and Zhou, Z.-H. Fast provably robust decision trees and boosting. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 8127–8144, Baltimore, MD, 2022.
- Huang, B., Ke, Z., Wang, Y., Wang, W., Shen, L., and Liu, F. Adversarial defence by diversified simultaneous training of deep ensembles. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, pp. 7823–7831, Virtual Event, 2021.
- Kariyappa, S. and Qureshi, M. K. Improving adversarial robustness of ensembles with diversity training. *CoRR*, abs/1901.09981, 2019.

- Katz, G., Barrett, C., Dill, D. L., Julian, K., and Kochenderfer, M. J. Reluplex: An efficient smt solver for verifying deep neural networks. In *Proceedings of the 29th International Conference on Computer Aided Verification*, pp. 97–117, Heidelberg, Germany, 2017.
- Krogh, A. and Vedelsby, J. Neural network ensembles, cross validation, and active learning. In *Advances in Neural Information Processing Systems 7*, pp. 231, Denver, CO, 1994.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *Proceedings of the 6th International Conference on Learning Representations*, Vancouver, Canada, 2018.
- Moré, J. J. and Sorensen, D. C. Computing a trust region step. *SIAM Journal on Scientific and Statistical Computing*, 4(3):553–572, 1983.
- Pang, T., Xu, K., Du, C., Chen, N., and Zhu, J. Improving adversarial robustness via promoting ensemble diversity. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 4970–4979, Long Beach, LA, 2019.
- Pang, T., Yang, X., Dong, Y., Su, H., and Zhu, J. Bag of tricks for adversarial training. In *Proceedings of the 9th International Conference on Learning Representations*, Virtual Event, 2021.
- Peng, S., Xu, W., Cornelius, C., Hull, M., Li, K., Duggal, R., Phute, M., Martin, J., and Chau, D. H. Robust principles: Architectural design principles for adversarially robust cnns. In *34th British Machine Vision Conference*, pp. 739–740, Aberdeen, UK, 2023.
- Pinot, R., Ettetdgui, R., Rizk, G., Chevaleyre, Y., and Atif, J. Randomization matters how to defend against strong adversarial attacks. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 7717–7727, Vienna, Austria, 2020.
- Robbins, H. and Monro, S. A stochastic approximation method. *Annals of Mathematical Statistics*, pp. 400–407, 1951.
- Rudin, W. et al. *Principles of Mathematical Analysis*, volume 3. McGraw-hill New York, 1964.
- Sen, S., Ravindran, B., and Raghunathan, A. Empir: Ensembles of mixed precision deep networks for increased robustness against adversarial attacks. In *Proceedings of the 8th International Conference on Learning Representations*, Addis Ababa, Ethiopia, 2020.
- Strauss, T., Hanselmann, M., Junginger, A., and Ulmer, H. Ensemble methods as a defense to adversarial perturbations against deep neural networks. *CoRR*, abs/1709.03423, 2017.
- Sun, T. and Zhou, Z.-H. Structural diversity of decision tree ensemble learning. *Frontiers of Computer Science*, 12(3): 560–570, 2018.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. In *Proceedings of the 2nd International Conference on Learning Representations*, Banff, Canada, 2014.
- Wang, Z., Pang, T., Du, C., Lin, M., Liu, W., and Yan, S. Better diffusion models further improve adversarial training. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 36246–36263, Honolulu, Hawaii, 2023.
- Weng, L., Zhang, H., Chen, H., Song, Z., Hsieh, C.-J., Daniel, L., Boning, D., and Dhillon, I. Towards fast computation of certified robustness for relu networks. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 5276–5285, Stockholm, Sweden, 2018.
- Wood, D., Mu, T., Webb, A. M., Reeve, H. W., Lujan, M., and Brown, G. A unified theory of diversity in ensemble learning. *Journal of Machine Learning Research*, 24 (359):1–49, 2023.
- Yang, H., Zhang, J., Dong, H., Inkawhich, N., Gardner, A., Touchet, A., Wilkes, W., Berry, H., and Li, H. Dverge: diversifying vulnerabilities for enhanced robust generation of ensembles. In *Advances in Neural Information Processing Systems 33*, pp. 5505–5515, Virtual Event, 2020.
- Yang, Z., Li, L., Xu, X., Zuo, S., Chen, Q., Zhou, P., Rubinstein, B., Zhang, C., and Li, B. Trs: Transferability reduced ensemble via promoting gradient diversity and model smoothness. In *Advances in Neural Information Processing Systems 34*, pp. 17642–17655, Virtual Event, 2021.
- Yang, Z., Li, L., Xu, X., Kailkhura, B., Xie, T., and Li, B. On the certified robustness for ensemble models and beyond. In *Proceedings of the 10th International Conference on Learning Representations*, Virtual Event, 2022.
- Young, L. An inequality of the hölder type, connected with stieltjes integration. *Acta Mathematica*, 67(1):251–282, 1936.

Zhang, D., Zhang, H., Courville, A., Bengio, Y., Ravikumar, P., and Suggala, A. S. Building robust ensembles via margin boosting. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 26669–26692, Baltimore, MD, 2022.

Zhang, H., Cheng, M., and Hsieh, C.-J. Enhancing certifiable robustness via a deep model ensemble. *CoRR*, abs/1910.14655, 2019a.

Zhang, H., Yu, Y., Jiao, J., Xing, E., El Ghaoui, L., and Jordan, M. Theoretically principled trade-off between robustness and accuracy. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 7472–7482, Long Beach, LA, 2019b.

Zhou, Z.-H. *Ensemble Methods: Foundations and Algorithms*. CRC Press, 2012.

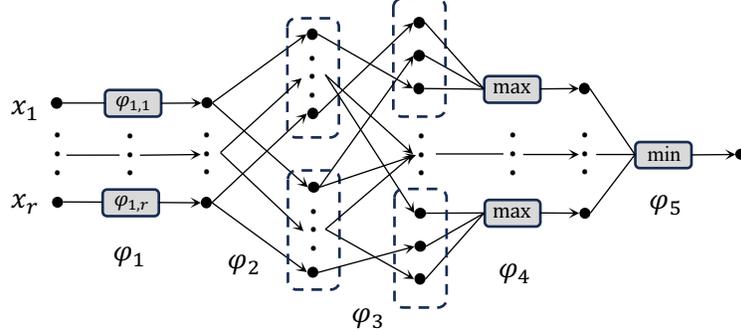


Figure 10. The structure of transformed neural network. The $\varphi_{1,j}(\cdot)$, $\max(\cdot)$ and $\min(\cdot)$ functions can be constructed by ReLU functions. The φ_2 is fully connected layer with $-1, 1$ or 0 weights. The φ_3 is fully connected layer with 1 or 0 weights.

A. Proof of Theorem 3.1

Definition A.1 (3-SAT problem). Given r boolean variables and s clauses in a conjunctive normal form CNF formula with each clause's size at most 3, is there an assignment to the r variables to make the CNF formula to be satisfied?

We present a key lemma for NP-hardness of adversarial squared loss, and the basic idea is a reduction from 3-SAT problem.

Lemma A.2. For some ReLU neural network $g(\cdot)$ and example (\mathbf{x}_0, y_0) , it is an NP-hard problem to precisely calculate the following adversarial squared loss

$$\max_{\delta \in \Delta_p^\epsilon} \{(g(\mathbf{x} + \delta) - y)^2\}.$$

Proof. It is sufficient to prove the NP-hardness of solving $\max_{\delta \in \Delta_p^\epsilon} \{(g(\mathbf{x} + \delta) - y)^2\} = \gamma$ for some $\gamma > 0$. Let $\phi = C_1 \wedge C_2 \wedge \dots \wedge C_s$ be a 3-SAT formula over a variable set $V = \{v_1, \dots, v_r\}$. Each $C_i = q_i^1 \vee q_i^2 \vee q_i^3$ is a disjunction with three literals q_i^1, q_i^2, q_i^3 , and each literal is a variable from V or their negations. The 3-SAT problem is to determine whether there exists an assignment $a : V \rightarrow \{0, 1\}$ for true ϕ .

We will show that any 3-SAT formula ϕ can be transformed into a neural network g over sample $(\mathbf{x}_0, y_0) = (\mathbf{0}, -1)$ in polynomial time, as well as the following sufficient and necessary condition:

$$\phi \text{ is satisfiable} \iff \max_{\delta \in \Delta_p^\epsilon} \{(g(\mathbf{x}_0 + \delta) - y_0)^2\} = (\epsilon/r^{1/p} + 1)^2. \quad (6)$$

We will construct the ReLU neural network $g(\mathbf{x}) = \varphi_5 \circ \varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x})$ with input $\mathbf{x} = (x_1, x_2, \dots, x_r) \in \mathbb{R}^r$ as shown in Figure 10. Specifically, we present the detailed constructions as follows:

- We construct function $\varphi_1(\mathbf{x}) : \mathbb{R}^r \rightarrow \mathbb{R}^r$ with the j -th element

$$[\varphi_1(\mathbf{x})]_j = \max(x_j + \epsilon/r^{1/p}, 0) - \max(x_j - \epsilon/r^{1/p}, 0) - \epsilon/r^{1/p} = \begin{cases} \epsilon/r^{1/p} & \text{for } x_j > \epsilon/r^{1/p} \\ x_j & \text{for } -\epsilon/r^{1/p} < x_j \leq \epsilon/r^{1/p} \\ -\epsilon/r^{1/p} & \text{for } x_j \leq -\epsilon/r^{1/p} \end{cases},$$

For $\mathbf{x}_0 = \mathbf{0}$ and $\|\delta\|_p \leq \epsilon$, we can achieve the maximum or minimum of the elements of $\varphi_1(\mathbf{x}_0 + \delta)$ independently. It is easy to construct $\varphi_1(\mathbf{x})$ with a ReLU neural network because of three operators $+$, $-$ and $\max(\cdot, 0)$. Intuitively, the i -th element of $\varphi_1(\mathbf{x}_0 + \delta)$ can be viewed as the variable v_i in 3-SAT problem, and its value $-\epsilon/r^{1/p}$ and $\epsilon/r^{1/p}$ can be viewed as the false and true of variable v_i , respectively.

- We construct $\varphi_2(\mathbf{x}) = (\mathbf{x}, -\mathbf{x})$, and this follows that, for $j \in [2r]$, $\mathbf{x}_0 = \mathbf{0}$ and $\|\delta\|_p \leq \epsilon$,

$$-\epsilon/r^{1/p} \leq [\varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)]_j \leq \epsilon/r^{1/p}.$$

We achieve the maximum or minimum for the first r elements of $\varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ independently. The last r elements are always the opposite of the first r elements. Intuitively, the last r elements of $\varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ can be viewed as the negation of variables v_1, \dots, v_r in 3-SAT problem.

- We construct $\varphi_3(\mathbf{x}) = W\mathbf{x}$ with $W = (\mathbf{w}_1^1; \mathbf{w}_1^2; \mathbf{w}_1^3; \dots; \mathbf{w}_s^1; \mathbf{w}_s^2; \mathbf{w}_s^3)^T \in \mathbb{R}^{3s \times 2r}$. Here, we construct three vectors $\mathbf{w}_i^1, \mathbf{w}_i^2, \mathbf{w}_i^3 \in \{0, 1\}^{2r}$ for clause $C_i = q_i^1 \vee q_i^2 \vee q_i^3$ for $i \in [s]$. For $k \in [3]$, the \mathbf{w}_i^k is the unit vector with j -th and $(j+r)$ -th element 1 if q_i^k is variable v_j and its negation $\neg v_j$, respectively. For $j \in [3s]$, $\mathbf{x}_0 = \mathbf{0}$ and $\|\delta\|_p \leq \epsilon$, we have

$$-\epsilon/r^{1/p} \leq [\varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)]_j \leq \epsilon/r^{1/p}.$$

Intuitively, every three elements of $\varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ can be viewed as three literals of a clause in 3-SAT problem. We achieve independently the maximum or minimum for $\varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ whose corresponding literals are different variables.

- We construct $\varphi_4(\mathbf{x}) = (\max\{x_1, x_2, x_3\}, \max\{x_4, x_5, x_6\}, \dots, \max\{x_{3s-2}, x_{3s-1}, x_{3s}\})$, where $\max\{\cdot, \cdot, \cdot\} = \max\{\max\{\cdot, \cdot\}, \cdot\}$ and $\min\{\cdot, \cdot, \cdot\} = \min\{\min\{\cdot, \cdot\}, \cdot\}$ can be constructed via ReLU function ($\max\{\cdot, 0\}$ function) as

$$\max\{a, b\} = \max\{a - b, 0\} + b \quad \text{and} \quad \min\{a, b\} = -\max\{b - a, 0\} + b. \quad (7)$$

For $j \in [s]$, $\mathbf{x}_0 = \mathbf{0}$ and $\|\delta\|_p \leq \epsilon$, we have

$$-\epsilon/r^{1/p} \leq [\varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)]_j \leq \epsilon/r^{1/p}.$$

We can achieve the maximum for the j -th element of $\varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ if and only if the $(3j-2)$ -th, $(3j-1)$ -th or $3j$ -th elements of $\varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ are maximal. Intuitively, the max function can be viewed as the \vee operator, and the i -th element of $\varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ can be viewed as the clause C_i in 3-SAT problem.

- We construct $\varphi_5(\mathbf{x}) = \min\{x_1, x_2, \dots, x_s\} = \min\{\min\{\min\{\dots, x_{s-2}\}, x_{s-1}\}, x_s\}$ via ReLU function. For $\mathbf{x}_0 = \mathbf{0}$ and $\|\delta\|_p \leq \epsilon$, we have

$$-\epsilon/r^{1/p} \leq \varphi_5 \circ \varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta) \leq \epsilon/r^{1/p}. \quad (8)$$

We achieve the maximum if and only if all elements of $\varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ are maximal. Intuitively, the min function can be viewed as the \wedge operator in 3-SAT problem. The $-\epsilon/r^{1/p}$ and $\epsilon/r^{1/p}$ value of $\varphi_5 \circ \varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\mathbf{x}_0 + \delta)$ can be viewed as the false and true of the CNF formula $\phi = C_1 \wedge C_2 \wedge \dots \wedge C_s$, respectively.

For $(\mathbf{x}_0, y_0) = (\mathbf{0}, -1)$ and $\epsilon \leq 1$, it remains to show that, from Eqns. (6) and (8),

$$\phi \text{ is satisfiable} \iff \text{there is an } \delta \text{ s.t. } g(\delta) = \epsilon/r^{1/p}. \quad (9)$$

\implies If $\phi = C_1 \wedge C_2 \wedge \dots \wedge C_m$ is satisfiable, then there is a satisfiable assignment α . We set the i -th element of δ as $\epsilon/r^{1/p}$ and $-\epsilon/r^{1/p}$ if $\alpha(v_i)$ is true and false, respectively. Then, we discuss $g(\delta)$ step by step as follows:

- The i -th element of $\varphi_1(\delta)$ is $-\epsilon/r^{1/p}$ and $\epsilon/r^{1/p}$ if the i -th element of δ is $-\epsilon/r^{1/p}$ and $\epsilon/r^{1/p}$, respectively;
- The first r elements of $\varphi_2 \circ \varphi_1(\delta)$ are equal to $\varphi_1(\delta)$ while the last r elements are the opposite of the first r elements;
- The $(3(i-1) + j)$ -th element of $\varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$ is the l -th and $(l+r)$ -th element of $\varphi_2 \circ \varphi_1(\delta)$ if literal q_i^j is variable v_l and $\neg v_l$, respectively;
- Every element of $\varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$ is $\epsilon/r^{1/p}$, since there is at least one true literal in q_i^1, q_i^2, q_i^3 for every satisfiable $C_i = q_i^1 \vee q_i^2 \vee q_i^3$, and hence there is at least one element with $\epsilon/r^{1/p}$ in every three elements of $\varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$;
- The final output $\varphi_5 \circ \varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$ is $\epsilon/r^{1/p}$, since φ_5 takes the minimum of $\varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$.

\impliedby If $g(\delta) = \epsilon/r^{1/p}$, we discuss δ as follows:

- Every elements of $\varphi_4 \circ \varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$ is $\epsilon/r^{1/p}$ from $\varphi_5(\cdot) \in [-\epsilon/r^{1/p}, \epsilon/r^{1/p}]$;
- At least one element is equal to $\epsilon/r^{1/p}$ in every three elements of $\varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$, since φ_4 takes the maximum of every three elements. Without loss of generality, let all $(3(i-1) + 1)$ -th elements be $\epsilon/r^{1/p}$ for $i \in [m]$;

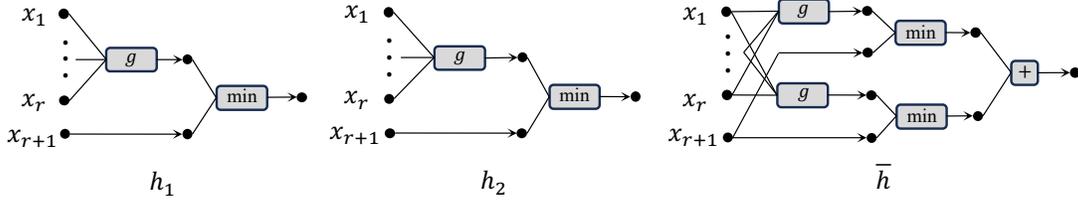


Figure 11. The structure of the transformed neural network. The φ block is a network similarly in the proof of Theorem A.2.

- For the first r elements of $\varphi_2 \circ \varphi_1(\delta)$, the l -th element will be equal to $\epsilon/r^{1/p}$ if literal q_i^1 is variable v_l . For the last r elements of $\varphi_2 \circ \varphi_1(\delta)$, the $(l+r)$ -th element will be equal to $\epsilon/r^{1/p}$ if literal q_i^1 is the negation of variable $\neg v_l$;
- The l -th element of $\varphi_1(\delta)$ is $\epsilon/r^{1/p}$ and $-\epsilon/r^{1/p}$ if the l -th element of $\varphi_2 \circ \varphi_1(\delta)$ is $\epsilon/r^{1/p}$ and $-\epsilon/r^{1/p}$, respectively;
- The l -th element of δ is $\epsilon/r^{1/p}$ and $-\epsilon/r^{1/p}$ if the l -th element of $\varphi_1(\delta)$ is $\epsilon/r^{1/p}$ and $-\epsilon/r^{1/p}$, respectively.

Let v_l be true and false if the l -th element of δ is $\epsilon/r^{1/p}$ and $-\epsilon/r^{1/p}$, respectively. This assignment makes the CNF formula satisfiable, since there is at least one element in every three elements of $\varphi_3 \circ \varphi_2 \circ \varphi_1(\delta)$ with $\epsilon/r^{1/p}$. \square

Proof of Theorem 3.1. It is sufficient to prove the NP-hardness of solving

$$\frac{1}{2} \sum_{j=1}^2 \max_{\delta \in \Delta_p^\epsilon} \{(f_j(\mathbf{x} + \delta) - y)^2\} - \max_{\delta \in \Delta_p^\epsilon} \{(\bar{f}(\mathbf{x} + \delta) - y)^2\} = \gamma \quad \text{for some } \gamma > 0.$$

Let $\phi = C_1 \wedge C_2 \wedge \dots \wedge C_s$ be a 3-SAT formula over variable set $V = \{v_1, \dots, v_r\}$. Each $C_i = q_i^1 \vee q_i^2 \vee q_i^3$ is a disjunction with three literals q_i^1, q_i^2, q_i^3 , and each literal is a variable from V or their negations. We will show that any 3-SAT formula ϕ can be transformed into two neural networks f_1, f_2 and sample $(x_0, y_0) = (\mathbf{0}, -1)$ in polynomial time, as well as the following sufficient and necessary condition:

$$\phi \text{ is satisfiable} \iff \Upsilon := \frac{1}{2} \sum_{j=1}^2 \max_{\delta \in \Delta_p^\epsilon} \{(f_j(\delta) + 1)^2\} - \max_{\delta \in \Delta_p^\epsilon} \{(\bar{f}(\delta) + 1)^2\} = (\epsilon/r^{1/p} + 1)^2 - 1. \quad (10)$$

We will construct the ReLU neural network $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ with inputs $\mathbf{x} \in \mathbb{R}^{r+1}$ in Figure 11. For $\phi = C_1 \wedge C_2 \wedge \dots \wedge C_s$, we construct neural network $g(\mathbf{x}_{1:r})$ as shown in Lemma A.2, where $\mathbf{x}_{1:r}$ are the first r elements of \mathbf{x} . We construct $f_1(\mathbf{x}) = \min\{g(\mathbf{x}_{1:r}), \eta(x_{r+1})\}$ with

$$\eta(x_{r+1}) = \max(x_{r+1} + \frac{\epsilon}{r^{1/p}}, 0) - \max(x_{r+1} - \frac{\epsilon}{r^{1/p}}, 0) - \frac{\epsilon}{r^{1/p}} = \begin{cases} \epsilon/r^{1/p} & \text{for } x_{r+1} > \epsilon/r^{1/p} \\ x_{r+1} & \text{for } -\epsilon/r^{1/p} < x_{r+1} \leq \epsilon/r^{1/p} \\ -\epsilon/r^{1/p} & \text{for } x_{r+1} \leq -\epsilon/r^{1/p}. \end{cases}$$

It is easy to construct $f_1(\mathbf{x})$ via a ReLU neural network from Eqn. (7), and we have

$$\phi \text{ is satisfiable} \iff \max_{\delta \in \Delta_p^\epsilon} \{(f_1(\delta) + 1)^2\} = (\epsilon/r^{1/p} + 1)^2. \quad (11)$$

This is because $-\epsilon/r^{1/p} \leq f_1(\mathbf{x}) \leq \epsilon/r^{1/p}$ from Eqn. (8), and for $\epsilon \in (0, 1]$, ϕ is satisfiable if and only if there exists an δ s.t. $f_1(\delta) = \epsilon/r^{1/p}$. If $f_1(\delta) = \epsilon/r^{1/p}$, then ϕ is satisfiable from Eqn. (9), since we have $g(\delta_{1:r}) = \epsilon/r^{1/p}$ from $f_1(\delta) = \min\{g(\delta_{1:r}), \eta(\delta_{r+1})\}$ and $\eta(\delta_{r+1}) \leq \epsilon/r^{1/p}$. For the inverse direction, if ϕ is satisfiable, then there is an $\delta_{1:r}$ s.t. $g(\delta_{1:r}) = \epsilon/r^{1/p}$ from Eqn. (9). This follows that $f_1(\delta) = \epsilon/r^{1/p}$ for $\delta_{k+1} = \epsilon/r^{1/p}$.

In a similar manner, we construct $f_2(\mathbf{x}) = \min\{g(\mathbf{x}_{1:k}), \eta(-x_{k+1})\}$, and have

$$\phi \text{ is satisfiable} \iff \max_{\delta \in \Delta_p^\epsilon} \{(f_2(\delta) + 1)^2\} = (\epsilon/r^{1/p} + 1)^2. \quad (12)$$

For odd function $\eta(\mathbf{x})$, we have the ensemble $\bar{f}(\mathbf{x}) = (f_1(\mathbf{x}) + f_2(\mathbf{x}))/2$ as

$$\bar{f}(\mathbf{x}) = \frac{1}{2}(\min\{g(\mathbf{x}_{1:k}), \eta(x_{k+1})\} + \min\{g(\mathbf{x}_{1:k}), \eta(-x_{k+1})\}) \leq \frac{1}{2}(\eta(x_{k+1}) + \eta(-x_{k+1})) = \frac{1}{2}(\eta(x_{k+1}) - \eta(x_{k+1})) = 0.$$

We have $-\epsilon/r^{1/p} \leq \bar{f}(\mathbf{x}) \leq 0$, and it holds that $(\bar{f}(\boldsymbol{\delta}) + 1)^2 \leq 1$ for $\epsilon \in (0, 1]$, and the equality holds for $\boldsymbol{\delta} = \mathbf{0}$. We have

$$\max_{\|\boldsymbol{\delta}\|_p \leq \epsilon} \{(\bar{f}(\boldsymbol{\delta}) + 1)^2\} = 1. \quad (13)$$

For the proof of Eqn. (10), we have $\Upsilon = (\epsilon/r^{1/p} + 1)^2 - 1$ if ϕ is satisfiable from Eqns. (11)-(13). For the inverse direction, if $\Upsilon = (\epsilon/r^{1/p} + 1)^2 - 1$, then we have, from Eqn. (13),

$$\frac{1}{2} \sum_{j=1}^2 \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \{(f_j(\boldsymbol{\delta}) + 1)^2\} = (\epsilon/r^{1/p} + 1)^2. \quad (14)$$

From $f_1(\boldsymbol{\delta}), f_2(\boldsymbol{\delta}) \in [-\epsilon/r^{1/p}, \epsilon/r^{1/p}]$ and $\epsilon \in (0, 1]$, we have

$$(f_1(\boldsymbol{\delta}) + 1)^2 \leq (\epsilon/r^{1/p} + 1)^2 \quad \text{and} \quad (f_2(\boldsymbol{\delta}) + 1)^2 \leq (\epsilon/r^{1/p} + 1)^2.$$

This follows that, from Eqn. (14)

$$\max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \{(f_1(\boldsymbol{\delta}) + 1)^2\} = (\epsilon/r^{1/p} + 1)^2 \quad \text{and} \quad \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \{(f_2(\boldsymbol{\delta}) + 1)^2\} = (\epsilon/r^{1/p} + 1)^2;$$

therefore, ϕ is satisfiable from Eqns. (11)-(12). This completes the proof. \square

B. Appendix for Section 3.2

B.1. Proof of Theorem 3.2

We begin with some useful lemmas as follows:

Lemma B.1 (Hölder's inequality (Young, 1936)). *For two real vectors $\mathbf{a} = (a_1, \dots, a_d)$ and $\mathbf{b} = (b_1, \dots, b_d)$, we have*

$$\sum_{i=1}^d |a_i b_i| \leq \|\mathbf{a}\|_p \|\mathbf{b}\|_q \quad \text{for positive } p \text{ and } q \text{ with } 1/p + 1/q = 1,$$

where the equality holds if and only if $\alpha|a_i|^p = \beta|b_i|^q$ for $i \in [d]$ w.r.t. some positive constants α and β .

Lemma B.2. *For vectors $\mathbf{w}, \boldsymbol{\delta} \in \mathbb{R}^d$, we have*

$$\max_{\|\boldsymbol{\delta}\|_p \leq \epsilon} \mathbf{w}^T \boldsymbol{\delta} = \epsilon \|\mathbf{w}\|_q \quad \text{and} \quad \min_{\|\boldsymbol{\delta}\|_p \leq \epsilon} \mathbf{w}^T \boldsymbol{\delta} = -\epsilon \|\mathbf{w}\|_q.$$

Proof. For every $\boldsymbol{\delta}$ with $\|\boldsymbol{\delta}\|_p \leq \epsilon$, we have, from Lemma B.1,

$$\mathbf{w}^T \boldsymbol{\delta} \leq \sum_{i=1}^d |w_i \delta_i| \leq \|\boldsymbol{\delta}\|_p \|\mathbf{w}\|_q \leq \epsilon \|\mathbf{w}\|_q.$$

Notice that the above equality holds if we choose $\boldsymbol{\delta} = \boldsymbol{\delta}^* = (\delta_1^*, \delta_2^*, \dots, \delta_d^*)$ with

$$\delta_i^* = \text{sign}(w_i) \epsilon (|w_i|^q / \|\mathbf{w}\|_q^q)^{1/p},$$

where $\text{sign}(w_i)$ is equal to $-1, 0, 1$ if w_i is negative, zero or positive, respectively. This follows that

$$\max_{\|\boldsymbol{\delta}\|_p \leq \epsilon} \mathbf{w}^T \boldsymbol{\delta} = \epsilon \|\mathbf{w}\|_q.$$

We also have, by letting $\boldsymbol{\delta}' = -\boldsymbol{\delta}$,

$$\min_{\|\boldsymbol{\delta}\|_p \leq \epsilon} \mathbf{w}^T \boldsymbol{\delta} = \min_{\|\boldsymbol{\delta}'\|_p \leq \epsilon} \mathbf{w}^T (-\boldsymbol{\delta}') = \min_{\|\boldsymbol{\delta}'\|_p \leq \epsilon} -\mathbf{w}^T \boldsymbol{\delta}' = -\max_{\|\boldsymbol{\delta}'\|_p \leq \epsilon} \mathbf{w}^T \boldsymbol{\delta}' = -\epsilon \|\mathbf{w}\|_q,$$

which completes the proof. \square

Lemma B.3. For linear function $h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, we have, for positive p, q with $1/p + 1/q = 1$,

$$\max_{\|\delta\|_p \leq \epsilon} (\mathbf{w}^T(\mathbf{x} + \delta) + b - y)^2 = (|\mathbf{w}^T \mathbf{x} + b - y| + \|\mathbf{w}\|_q \epsilon)^2.$$

Proof. We have the adversarial squared loss

$$\max_{\|\delta\|_p \leq \epsilon} (\mathbf{w}^T(\mathbf{x} + \delta) + b - y)^2 = \max_{\|\delta\|_p \leq \epsilon} (\mathbf{w}^T \mathbf{x} + b - y + \mathbf{w}^T \delta)^2.$$

From Lemma B.2, we have $-\|\mathbf{w}\|_q \epsilon \leq \mathbf{w}^T \delta \leq \|\mathbf{w}\|_q \epsilon$ and

$$\max_{\|\delta\|_p \leq \epsilon} (\mathbf{w}^T(\mathbf{x} + \delta) + b - y)^2 = \max_{\|\delta\|_p \leq \epsilon} (\mathbf{w}^T \mathbf{x} + b - y + \mathbf{w}^T \delta)^2 = (|\mathbf{w}^T \mathbf{x} + b - y| + \|\mathbf{w}\|_q \epsilon)^2,$$

which completes the proof. \square

Proof of Theorem 3.2. We have the adversarial loss for f_j and \bar{f} , from Lemma B.3,

$$\begin{aligned} \max_{\delta \in \Delta_p^\epsilon} \{(\bar{f}(\mathbf{x} + \delta) - y)^2\} &= \max_{\delta \in \Delta_p^\epsilon} \{(\bar{f}(\mathbf{x}) + \nabla \bar{f}(\mathbf{x})^T \delta - y)^2\} = (|\bar{f}(\mathbf{x}) - y| + \|\nabla \bar{f}(\mathbf{x})\|_q \epsilon)^2, \\ \max_{\delta \in \Delta_p^\epsilon} \{(f_j(\mathbf{x} + \delta) - y)^2\} &= \max_{\delta \in \Delta_p^\epsilon} \{(f_j(\mathbf{x}) + \nabla f_j(\mathbf{x})^T \delta - y)^2\} = (|f_j(\mathbf{x}) - y| + \|\nabla f_j(\mathbf{x})\|_q \epsilon)^2. \end{aligned}$$

This follows that

$$\begin{aligned} &\frac{1}{m} \sum_{j=1}^m \max_{\delta \in \Delta_p^\epsilon} \{(f_j(\mathbf{x} + \delta) - y)^2\} - \max_{\delta \in \Delta_p^\epsilon} \{(\bar{f}(\mathbf{x} + \delta) - y)^2\} \\ &= \frac{1}{m} \sum_{j=1}^m (|f_j(\mathbf{x}) - y| + \|\nabla f_j(\mathbf{x})\|_q \epsilon)^2 - (|\bar{f}(\mathbf{x}) - y| + \|\nabla \bar{f}(\mathbf{x})\|_q \epsilon)^2 \\ &\quad - \left((\bar{f}(\mathbf{x}) - y)^2 + 2|\bar{f}(\mathbf{x}) - y| \|\nabla \bar{f}(\mathbf{x})\|_q \epsilon + \|\nabla \bar{f}(\mathbf{x})\|_q^2 \epsilon^2 \right) \\ &= \frac{\epsilon^2}{m} \sum_{j=1}^m (\|\nabla f_j(\mathbf{x})\|_q^2 - \|\nabla \bar{f}(\mathbf{x})\|_q^2) + \frac{2\epsilon}{m} \sum_{j=1}^m (\|\nabla f_j(\mathbf{x})\|_q |f_j(\mathbf{x}) - y| - \|\nabla \bar{f}(\mathbf{x})\|_q |\bar{f}(\mathbf{x}) - y|) \\ &\quad + \frac{1}{m} \sum_{j=1}^m (f_j(\mathbf{x}) - y)^2 - (\bar{f}(\mathbf{x}) - y)^2 \\ &= \mathbf{Gradient Diversity} + \mathbf{Cross Diversity} + \frac{1}{m} \sum_{j=1}^m (f_j(\mathbf{x}) - y)^2 - (\bar{f}(\mathbf{x}) - y)^2. \end{aligned}$$

We also have

$$\frac{1}{m} \sum_{j=1}^m (f_j(\mathbf{x}) - y)^2 - (\bar{f}(\mathbf{x}) - y)^2 = \frac{1}{m} \sum_{j=1}^m f_j(\mathbf{x})^2 - \bar{f}(\mathbf{x})^2 = \frac{1}{m} \sum_{j=1}^m (f_j(\mathbf{x}) - \bar{f}(\mathbf{x}))^2,$$

which completes the proof. \square

B.2. Training Details for Diversity Decomposition

We consider the base learners as convolutional neural network with two convolutional layers and one MLP layer of 100 neurons. The first convolutional layer has 24 filters with kernel size 5, while the second convolutional layer has 24 filters of kernel size 5. We take the ReLU activation function, and the input and output sizes are $3 \times 32 \times 32$ and 1, respectively.

We select the 5-th and 6-th class on datasets MNIST and F-MNIST to train the base learners independently, and take the SGD method (Robbins & Monro, 1951) with batch size 256 and learning rate 0.01. We consider the PGD-attack (Madry et al., 2018) to calculate the adversarial ensemble loss and average of individual adversarial losses. The perturbation size is set to $8/255$ and $128/255$ for l_∞ and l_2 norm, respectively.

B.3. Discussions on the Second-Order Approximation

For the second-order approximation, we have

$$f(\mathbf{x} + \boldsymbol{\delta}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T H(\mathbf{x}) \boldsymbol{\delta},$$

where $H(\mathbf{x})$ is the Hessian matrix of f at point \mathbf{x} . This follows that

$$\max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} (f(\mathbf{x} + \boldsymbol{\delta}) - y)^2 = \max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} \left(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T H(\mathbf{x}) \boldsymbol{\delta} - y \right)^2,$$

and hence, it is sufficient to consider two problems as follows

$$\max_{\boldsymbol{\delta} \in \Delta_p^\epsilon} f(\mathbf{x}) + \nabla f(\mathbf{x})^T \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T H(\mathbf{x}) \boldsymbol{\delta} - y \quad \text{and} \quad \min_{\boldsymbol{\delta} \in \Delta_p^\epsilon} f(\mathbf{x}) + \nabla f(\mathbf{x})^T \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^T H(\mathbf{x}) \boldsymbol{\delta} - y.$$

We focus on the special case $p = 2$, and the above two problems can be formalized as

$$\min_{\|\boldsymbol{\delta}\|_2 \leq \epsilon} \frac{1}{2} \boldsymbol{\delta}^T B \boldsymbol{\delta} - \mathbf{b}^T \boldsymbol{\delta} \quad \text{for some } \mathbf{b} \in \mathbb{R}^n \text{ and symmetric matrix } B \in \mathbb{R}^{n \times n}. \quad (15)$$

This is known as the trust region subproblem (Forsythe & Golub, 1965; Moré & Sorensen, 1983; Fortin & Wolkowicz, 2004), and we have

- the optimal solution $\boldsymbol{\delta}^* = B^{-1} \mathbf{b}$ if there is no solution on the boundary $\{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_2 \leq \epsilon\}$ in (15), from the positive definiteness of B and $\|B^{-1} \mathbf{b}\| < \epsilon$;
- the optimal solution $\boldsymbol{\delta}^* = -(B + \alpha^* I)^{-1} \mathbf{b}$ if there is solution on the boundary $\{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_2 \leq \epsilon\}$ in (15) and $\alpha^* > -\lambda_1$. Here, λ_1 is the smallest eigenvalue of B and α^* is the solution of

$$\sum_{j=1}^n \frac{\gamma_j^2}{(\lambda_j + \alpha^*)^2} = \epsilon^2, \quad (16)$$

where $\lambda_1, \dots, \lambda_n$ are eigenvalues of B and $\gamma_1, \dots, \gamma_n$ are the elements of $Q^T \mathbf{b}$ with eigendecomposition $B = Q \Lambda Q^T$.

- the optimal solution $\boldsymbol{\delta}^* = \boldsymbol{\delta}_0 + \tau \mathbf{z}$ if there is solution on the boundary $\{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_2 \leq \epsilon\}$ in (15) and $\alpha^* \leq -\lambda_1$ in (16). Here, $\boldsymbol{\delta}_0$ is the solution of

$$(B - \lambda_1 I) \boldsymbol{\delta}_0 = -\mathbf{b} \quad \text{s.t.} \quad \|\boldsymbol{\delta}_0\| \leq \epsilon,$$

and \mathbf{z} is an eigenvector of B with eigenvalue λ_1 and $\tau \in \mathbb{R}$ satisfies $\|\boldsymbol{\delta}_0 + \tau \mathbf{z}\| = \epsilon$.

Here, the main challenge is how to obtain the closed-form solution of adversarial loss via second-order approximation, since it is relevant to the roots of high-order polynomials Eqn. (16).

B.4. Discussions of Average of \cos Values for l_∞ Norm

We focus on 2-dimensional instance space $\mathcal{X} \subseteq \mathbb{R}^2$ and label space $\mathcal{Y} \subseteq \mathbb{R}$, and consider

$$f_1(\mathbf{x}) = x_1 + x_2, \quad f_2(\mathbf{x}) = x_1 - 3x_2, \quad f_3(\mathbf{x}) = ax_1 + bx_2 \quad \text{and} \quad f_4(\mathbf{x}) = ax_1 - bx_2,$$

where $a = (\sqrt{5} - 1)/2$ and $b = \sqrt{2} + \sqrt{5}/5$. We study two ensembles: one ensemble of f_1 and f_2 ; the other ensemble of f_3 and f_4 . For example $(\mathbf{x}, y) = ([1, 0], 1)$ and perturbation set $\Delta = \{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_\infty \leq 1\}$, we have

$$\cos(\nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x})) = \frac{\langle \nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x}) \rangle}{\|\nabla f_1(\mathbf{x})\|_2 \|\nabla f_2(\mathbf{x})\|_2} = \frac{\langle \nabla f_3(\mathbf{x}), \nabla f_4(\mathbf{x}) \rangle}{\|\nabla f_3(\mathbf{x})\|_2 \|\nabla f_4(\mathbf{x})\|_2} = \cos(\nabla f_3(\mathbf{x}), \nabla f_4(\mathbf{x})),$$

and we also have the same average of individual adversarial losses from Lemma B.3 as follows

$$\sum_{i=1}^2 \max_{\boldsymbol{\delta} \in \Delta} \frac{(f_i(\mathbf{x} + \boldsymbol{\delta}) - y)^2}{2} = \sum_{i=3}^4 \max_{\boldsymbol{\delta} \in \Delta} \frac{(f_i(\mathbf{x} + \boldsymbol{\delta}) - y)^2}{2}.$$

However, the adversarial ensemble losses are different from

$$\max_{\delta \in \Delta} (f_1(\mathbf{x} + \delta)/2 + f_2(\mathbf{x} + \delta)/2 - y)^2 = 8 \quad \text{and} \quad \max_{\delta \in \Delta} (f_3(\mathbf{x} + \delta)/2 + f_4(\mathbf{x} + \delta)/2 - y)^2 \approx 7.2 .$$

Thus, there exist two ensembles of the same averages of cos values and individual adversarial losses, but with different adversarial ensemble losses for l_∞ norm.

C. Appendix for Section 3.3

C.1. Proof of Lemma 3.3

We have the cross-entropy loss

$$\ell(f(\mathbf{x} + \delta), y) = -y(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta) + \ln(1 + \exp(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta)) ,$$

and the adversarial cross-entropy loss

$$\max_{\delta \in \Delta_p^\epsilon} \ell(f(\mathbf{x} + \delta), y) = \max_{\delta \in \Delta_p^\epsilon} \{ -y(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta) + \ln(1 + \exp(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta)) \} . \quad (17)$$

For $\delta \in \Delta_p^\epsilon$, we have, from Lemma B.1

$$-\|\nabla f(\mathbf{x})\|_{q\epsilon} \leq \nabla f(\mathbf{x})^T \delta \leq \|\nabla f(\mathbf{x})\|_{q\epsilon} .$$

We get the maximum of Eqn. (17) when

$$\nabla f(\mathbf{x})^T \delta = \begin{cases} \|\nabla f(\mathbf{x})\|_{q\epsilon} & \text{for } y = 0 \\ -\|\nabla f(\mathbf{x})\|_{q\epsilon} & \text{for } y = 1 , \end{cases}$$

and the optimal adversarial perturbation δ^* with $\nabla f(\mathbf{x})^T \delta^* = -(2y - 1)\|\nabla f(\mathbf{x})\|_{q\epsilon}$. We finally have

$$f(\mathbf{x} + \delta^*) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta^* = f(\mathbf{x}) - (2y - 1)\|\nabla f(\mathbf{x})\|_{q\epsilon} ,$$

and the probability of the positive class of the adversarial example $\mathbf{x} + \delta^*$

$$p_{f,+}^{adv} = \frac{1}{1 + \exp(-(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \delta^*))} = \frac{1}{1 + \exp(-(f(\mathbf{x}) - (2y - 1)\|\nabla f(\mathbf{x})\|_{q\epsilon}))} ,$$

which completes the proof. \square

C.2. Proof of Theorem 3.4

We have

$$\begin{aligned} & \ln(1 + \exp(f_j(\mathbf{x}) - y'\|\nabla f_j(\mathbf{x})\|_{q\epsilon})) - \ln(1 + \exp(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon})) \\ &= \frac{1}{1 + \exp(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon})} \ln\left(\frac{1 + \exp(f_j(\mathbf{x}) - y'\|\nabla f_j(\mathbf{x})\|_{q\epsilon})}{1 + \exp(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon})}\right) \\ & \quad + \frac{1}{1 + \exp(-(f(\mathbf{x}) - y'\|\nabla f(\mathbf{x})\|_{q\epsilon}))} \ln\left(\frac{1 + \exp(f_j(\mathbf{x}) - y'\|\nabla f_j(\mathbf{x})\|_{q\epsilon})}{1 + \exp(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon})}\right) \\ &= \frac{1}{1 + \exp(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon})} \ln\left(\frac{1 + \exp(f_j(\mathbf{x}) - y'\|\nabla f_j(\mathbf{x})\|_{q\epsilon})}{1 + \exp(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon})}\right) \\ & \quad + \frac{1}{1 + \exp(-(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon}))} \ln\left(\frac{1 + \exp(-(f_j(\mathbf{x}) - y'\|\nabla f_j(\mathbf{x})\|_{q\epsilon}))}{1 + \exp(-(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon}))}\right) \\ & \quad + \frac{1}{1 + \exp(-(\bar{f}(\mathbf{x}) - y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon}))} (f_j(\mathbf{x}) - y'\|\nabla f_j(\mathbf{x})\|_{q\epsilon} - \bar{f}(\mathbf{x}) + y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon}) \\ &= KL(p_{\bar{f},adv}, p_{f_j,adv}) + p_{\bar{f},+,adv} (f_j(\mathbf{x}) - y'\|\nabla f_j(\mathbf{x})\|_{q\epsilon} - \bar{f}(\mathbf{x}) + y'\|\nabla \bar{f}(\mathbf{x})\|_{q\epsilon}) . \end{aligned}$$

This follows that, from Lemma 3.3,

$$\begin{aligned}
 & \frac{1}{m} \sum_{j=1}^m \max_{\|\delta\|_p \leq \epsilon} \ell(\tilde{f}_j(\mathbf{x} + \delta), y) - \max_{\|\delta\|_p \leq \epsilon} \ell(\bar{g}(\mathbf{x} + \delta), y) \\
 &= \frac{1}{m} \sum_{j=1}^m \ln(1 + \exp(f_j(\mathbf{x}) - y' \|\nabla f_j(\mathbf{x})\|_q \epsilon)) - y(f_j(\mathbf{x}) - y' \|\nabla f_j(\mathbf{x})\|_q \epsilon) \\
 &\quad - \ln(1 + \exp(\bar{f}(\mathbf{x}) - y' \|\nabla \bar{f}(\mathbf{x})\|_q \epsilon)) + y(\bar{f}(\mathbf{x}) - y' \|\nabla \bar{f}(\mathbf{x})\|_q \epsilon) \\
 &= \frac{1}{m} \sum_{j=1}^m KL(p_{\bar{f},adv}, p_{f_j,adv}) + (p_{\bar{f},+,adv} - y)(2y - 1) \epsilon \frac{1}{m} \sum_{j=1}^m (\|\nabla \bar{f}(\mathbf{x})\|_q - \|\nabla f_j(\mathbf{x})\|_q),
 \end{aligned}$$

which completes the proof. \square

C.3. Discussions of Average of cos Values for Cross-Entropy Loss

Example 2. There exist two ensembles of the same averages of cos values and individual adversarial cross-entropy losses, but with different adversarial ensemble losses for cross-entropy loss and l_2 norm.

Proof. We focus on 2-dimensional instance space $\mathcal{X} \subseteq \mathbb{R}^2$ and label space $\mathcal{Y} \subseteq \mathbb{R}$, and consider

$$f_1(\mathbf{x}) = x_1 + x_2, \quad f_2(\mathbf{x}) = x_1 - 3x_2, \quad f_3(\mathbf{x}) = abx_1 + bx_2 \quad \text{and} \quad f_4(\mathbf{x}) = abx_1 - bx_2,$$

where

$$a = \frac{\sqrt{5} - 1}{2} \quad \text{and} \quad b = \frac{2 \ln \left(\sqrt{(1 + \exp(1 + \sqrt{2}))(1 + \exp(1 + \sqrt{10}))} - 1 \right)}{\sqrt{5} - 1 + \sqrt{10} - 2\sqrt{5}}.$$

We study two ensembles: one ensemble of f_1 and f_2 ; the other ensemble of f_3 and f_4 . For example $(\mathbf{x}, y) = ([1, 0], 1)$ and perturbation set $\Delta = \{\delta : \|\delta\|_2 \leq 1\}$, we have

$$\cos(\nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x})) = \frac{\langle \nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x}) \rangle}{\|\nabla f_1(\mathbf{x})\|_2 \|\nabla f_2(\mathbf{x})\|_2} = \frac{\langle \nabla f_3(\mathbf{x}), \nabla f_4(\mathbf{x}) \rangle}{\|\nabla f_3(\mathbf{x})\|_2 \|\nabla f_4(\mathbf{x})\|_2} = \cos(\nabla f_3(\mathbf{x}), \nabla f_4(\mathbf{x})),$$

and we also have the same average of individual adversarial losses from Lemma 3.3

$$\sum_{i=1}^2 \max_{\delta \in \Delta} \frac{(f_i(\mathbf{x} + \delta) - y)^2}{2} = \sum_{i=3}^4 \max_{\delta \in \Delta} \frac{(f_i(\mathbf{x} + \delta) - y)^2}{2} = \frac{\ln(1 + \exp(1 + \sqrt{2})) + \ln(1 + \exp(1 + \sqrt{10}))}{2}.$$

However, the adversarial ensemble losses are different from

$$\max_{\delta \in \Delta} (f_1(\mathbf{x} + \delta)/2 + f_2(\mathbf{x} + \delta)/2 - y)^2 \approx 2.4999 \quad \text{and} \quad \max_{\delta \in \Delta} (f_3(\mathbf{x} + \delta)/2 + f_4(\mathbf{x} + \delta)/2 - y)^2 \approx 2.3738,$$

which completes the proof. \square

Example 3. There exist two ensembles of the same averages of cos values and individual adversarial cross-entropy losses, but with different adversarial ensemble losses for cross-entropy loss for l_∞ norm.

Proof. We focus on 2-dimensional instance space $\mathcal{X} \subseteq \mathbb{R}^2$ and label space $\mathcal{Y} \subseteq \mathbb{R}$, and consider

$$f_1(\mathbf{x}) = x_1 + x_2, \quad f_2(\mathbf{x}) = x_1 - 3x_2, \quad f_3(\mathbf{x}) = abx_1 + bx_2 \quad \text{and} \quad f_4(\mathbf{x}) = abx_1 - bx_2,$$

where $a = (\sqrt{5} - 1)/2$ and $b = \ln(\sqrt{(1 + \exp(3))(1 + \exp(5))} - 1)/\sqrt{5}$. We study two ensembles: one ensemble of f_1 and f_2 ; the other of f_3 and f_4 . For example $(\mathbf{x}, y) = ([1, 0], 1)$ and perturbation set $\Delta = \{\delta : \|\delta\|_\infty \leq 1\}$, we have

$$\cos(\nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x})) = \frac{\langle \nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x}) \rangle}{\|\nabla f_1(\mathbf{x})\|_2 \|\nabla f_2(\mathbf{x})\|_2} = \frac{\langle \nabla f_3(\mathbf{x}), \nabla f_4(\mathbf{x}) \rangle}{\|\nabla f_3(\mathbf{x})\|_2 \|\nabla f_4(\mathbf{x})\|_2} = \cos(\nabla f_3(\mathbf{x}), \nabla f_4(\mathbf{x})),$$

and we also have the same average of individual adversarial losses from Lemma 3.3

$$\sum_{i=1}^2 \max_{\delta \in \Delta} \frac{(f_i(\mathbf{x} + \delta) - y)^2}{2} = \sum_{i=3}^4 \max_{\delta \in \Delta} \frac{(f_i(\mathbf{x} + \delta) - y)^2}{2} = \frac{\ln(1 + \exp(3)) + \ln(1 + \exp(5))}{2}.$$

However, the adversarial ensemble losses are different from

$$\max_{\delta \in \Delta} (f_1(\mathbf{x} + \delta)/2 + f_2(\mathbf{x} + \delta)/2 - y)^2 \approx 3.0486 \quad \text{and} \quad \max_{\delta \in \Delta} (f_3(\mathbf{x} + \delta)/2 + f_4(\mathbf{x} + \delta)/2 - y)^2 \approx 2.3738,$$

which completes the proof. \square

D. Appendix for Section 4

D.1. Proof of the Extended Diversity

For simplicity, we abbreviate $\mathbf{f}(\tilde{\mathbf{x}}_{\mathbf{f}})$ and $\mathbf{p}_{\mathbf{f}}(\tilde{\mathbf{x}}_{\mathbf{f}})$ to \mathbf{f} and $\mathbf{p}_{\mathbf{f}}$, respectively. Let $f_k, p_{f_k}, k \in [K]$ be the k -th element of \mathbf{f} and $\mathbf{p}_{\mathbf{f}}$, respectively. We have the adversarial cross-entropy loss for multi-classification

$$\ell(\mathbf{f}(\tilde{\mathbf{x}}_{\mathbf{f}}), y) = -f_y + \log \left(\sum_{k=1}^K \exp(f_k) \right).$$

This follows that

$$\ell(\bar{\mathbf{f}}(\tilde{\mathbf{x}}_{\bar{\mathbf{f}}}), y) - \frac{1}{m} \sum_{j=1}^m \ell(\mathbf{f}_j(\tilde{\mathbf{x}}_{\mathbf{f}_j}), y) = -\bar{f}_y + \frac{1}{m} \sum_{j=1}^m f_{j,y} + \log \left(\sum_{k=1}^K \exp(\bar{f}_k) \right) - \frac{1}{m} \sum_{j=1}^m \log \left(\sum_{k=1}^K \exp(f_{j,k}) \right). \quad (18)$$

We also have

$$\begin{aligned} \log \left(\sum_{k=1}^K \exp(\bar{f}_k) \right) - \log \left(\sum_{k=1}^K \exp(f_{j,k}) \right) &= \sum_{k_1=1}^K \frac{\exp(\bar{f}_{k_1})}{\sum_{k=1}^K \exp(\bar{f}_k)} \log \left(\frac{\sum_{k=1}^K \exp(\bar{f}_k)}{\sum_{k=1}^K \exp(f_{j,k})} \right) \\ &= \sum_{k_1=1}^K \frac{\exp(\bar{f}_{k_1})}{\sum_{k=1}^K \exp(\bar{f}_k)} \log \left(\frac{f_{j,k_1} / \sum_{k=1}^K \exp(f_{j,k})}{\bar{f}_{k_1} / \sum_{k=1}^K \exp(\bar{f}_k)} \times \frac{\bar{f}_{k_1}}{f_{j,k_1}} \right) = \sum_{k=1}^K p_{\bar{f}_k} (\bar{f}_k - f_{j,k}) - \text{KL}(\mathbf{p}_{\bar{\mathbf{f}}}, \mathbf{p}_{\mathbf{f}_j}), \end{aligned}$$

and this follows that, by setting $\mathbf{r} = \mathbf{p}_{\bar{\mathbf{f}}}(\tilde{\mathbf{x}}_{\bar{\mathbf{f}}}) - \mathbf{e}_y$,

$$\frac{1}{m} \sum_{j=1}^m \log \left(\sum_{k=1}^K \exp(f_{j,k}) \right) - \log \left(\sum_{k=1}^K \exp(\bar{f}_k) \right) = \sum_{j=1}^m \frac{\langle \mathbf{r}, \mathbf{f}_j(\tilde{\mathbf{x}}_{\mathbf{f}_j}) - \bar{\mathbf{f}}(\tilde{\mathbf{x}}_{\bar{\mathbf{f}}}) \rangle}{m} + \sum_{j=1}^m \frac{\text{KL}(\mathbf{p}_{\bar{\mathbf{f}}}(\tilde{\mathbf{x}}_{\bar{\mathbf{f}}}), \mathbf{p}_{\mathbf{f}_j}(\tilde{\mathbf{x}}_{\mathbf{f}_j}))}{m}. \quad (19)$$

We have, from the first-order approximation $\mathbf{f}(\tilde{\mathbf{x}}_{\mathbf{f}}) \approx \mathbf{f}(\mathbf{x}) + \mathbf{J}_{\mathbf{f}}(\mathbf{x})\delta_{\mathbf{f}}^*$,

$$\frac{1}{m} \sum_{j=1}^m \langle \mathbf{r}, \mathbf{f}_j(\tilde{\mathbf{x}}_{\mathbf{f}_j}) - \bar{\mathbf{f}}(\tilde{\mathbf{x}}_{\bar{\mathbf{f}}}) \rangle = \frac{1}{m} \sum_{j=1}^m \langle \mathbf{r}, \mathbf{J}_{\mathbf{f}_j}(\mathbf{x})\delta_{\mathbf{f}_j}^* - \mathbf{J}_{\bar{\mathbf{f}}}(\mathbf{x})\delta_{\bar{\mathbf{f}}}^* \rangle,$$

which completes the proof by combining with Eqns. (18)-(19). \square

D.2. Proof of Orthogonality

We now show the orthogonalization of the predictions of base learners by optimizing Eqn. (5) as follows.

Theorem D.1. *For K multi-class learning, let f_1, \dots, f_m be m base learners with cross-entropy loss being at least B , and $K - 1$ is a multiple of m . We have the minimizer of Eqn. (5) over example (\mathbf{x}, y) as*

$$p_{f_i}(\mathbf{x} + \delta_{f_i}^*)_k = \begin{cases} \exp(-B) & \text{for } k = y \\ \frac{1 - \exp(-B)}{K-1} & \text{for } k \in s_i \\ 0 & \text{otherwise,} \end{cases} \quad (20)$$

where s_1, \dots, s_m is a partition of set $\{1, \dots, K\} \setminus \{y\}$ and $p_{f_i}(\tilde{\mathbf{x}}_{f_i})_k$ is the k -th element of $p_{f_i}(\tilde{\mathbf{x}}_{f_i})$.

Table 4. Comparison of training time (times(s) per epoch) for our AdvE_{OAP} with and without the regularization Γ .

Our AdvE _{OAP}	MNIST	F-MNIST	CIFAR10
without regularization $\Gamma_\alpha(\cdot)$	40.08	40.11	123.84
with regularization $\Gamma_\alpha(\cdot)$	41.45	41.32	125.48

Proof. We first have $p_{f_i}(\mathbf{x} + \boldsymbol{\delta}_{f_i}^*)_y = \exp(-B)$, since the adversarial cross-entropy loss is at least B for each base learner. The regularization in Eqn. (5) is defined as

$$\Gamma_\alpha(\mathbf{x}, y) = H\left(\sum_{j=1}^m \tilde{\mathbf{p}}_{f_j}(\tilde{\mathbf{x}}_{f_j})/m\right) + \alpha \cdot \log(V(\tilde{\mathbf{p}}_{f_1}(\tilde{\mathbf{x}}_{f_1}), \dots, \tilde{\mathbf{p}}_{f_m}(\tilde{\mathbf{x}}_{f_m}))).$$

The $V(\tilde{\mathbf{p}}_{f_1}(\tilde{\mathbf{x}}_{f_1}), \dots, \tilde{\mathbf{p}}_{f_m}(\tilde{\mathbf{x}}_{f_m}))$ achieves its maximum if and only if the non-label probability vectors of each individual network are mutually orthogonal (Bernstein, 2009). The $H(\sum_{j=1}^m \tilde{\mathbf{p}}_{f_j}(\tilde{\mathbf{x}}_{f_j})/m)$ achieves the maximum if and only if the mean of non-label probability vectors of individual networks are uniform. It is obvious that Eqn. (20) satisfies the two conditions simultaneously. Thus, Eqn. (20) is the minimizer of Eqn. (5). \square

D.3. Other Details of Algorithm 1

PGD-attack for l_p norm perturbation ball

The PGD-attack generates adversarial examples iteratively for l_∞ -norm perturbation ball as follows

$$\mathbf{x}^{t+1,j} = \prod_{\mathbf{x} + \Delta_\infty^\epsilon} (\mathbf{x}^{t,j} + \alpha \cdot \text{sign}(\nabla \ell(f_j(\mathbf{x}^{t,j}), y))),$$

where \prod denotes the projection and $\mathbf{x} + \Delta_\infty^\epsilon = \{\mathbf{x} + \boldsymbol{\delta} \mid \boldsymbol{\delta} \in \Delta_\infty^\epsilon\}$. The $\mathbf{x}^{0,j}$ is initialized as \mathbf{x} , and $\mathbf{x}^{T,j}$ is used as the adversarial example of base learner f_j . For other l_p norm perturbation ball, PGD-attack generates adversarial examples as

$$\mathbf{x}^{t+1,j} = \prod_{\mathbf{x} + \Delta_p^\epsilon} (\mathbf{x}^{t,j} + \alpha \cdot g_p(\nabla \ell(f_j(\mathbf{x}^{t,j}), y))),$$

where g_p is the function that maps the gradient to the update direction

$$g_p(\mathbf{w}) = \text{sign}(w_i) \epsilon (|w_i|^q / \|\mathbf{w}\|_q^q)^{1/p}.$$

Calculating the volume of polytope

For m vectors $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d$ and $X = (\mathbf{x}_1, \dots, \mathbf{x}_m) \in \mathbb{R}^{d \times m}$, we have, from matrix theory (Bernstein, 2009),

$$V^2(\mathbf{x}_1, \dots, \mathbf{x}_m) = \det(X^T X),$$

where $\det(X^T X)$ is the determination of the matrix $X^T X$. For matrix $A \in \mathbb{R}^{n \times n}$, we also have

$$\frac{\partial \det(A)}{\partial a_{ij}} = \det(A) (A^{-1})^T,$$

where a_{ij} is the i -th row and j -th column element of A . We finally optimize the objective Eqn. (5) for neural networks with SGD method (Robbins & Monro, 1951).

Time Complexity of Algorithm 1

The time complexity of Algorithm 1 takes m -times as that of training a single neural network adversarially (m is the number of neural networks in the ensemble). In addition, it takes $O(m^3)$ computational cost for the regularization with its gradient. In practice, the regularization takes much smaller computational cost than that of training neural networks, as in Table 4.

Table 5. Hyperparameters of all ensemble methods used in our experiments. Parameters that were not applicable were left blank.

Parameter	GAL	ADP	AdvADP	DVERGE	PDD	TRS	iGAT(ADP)
α	0.5	2	2	-	0.01	1	2
β	-	0.5	0.5	-	-	5	0.5

E. Appendix for Section 5

E.1. Experimental settings

For $iGAT_{ADP}$, we take 150, 150 and 480 epoches for MNIST, F-MNIST and CIFAR10 for convergence; while for other ensemble methods, we take 60, 60 and 250 epoches for MNIST, F-MNIST and CIFAR10, respectively. For adversarial examples in training process, we take PGD10 with 10 steps with step-size 0.04, 0.01 and 0.008 for MNIST, F-MNIST and CIFAR10, respectively. We set $\alpha = 0.02$ and $\lambda = 10$ for our method, and Table 5 summarizes parameter setting for others.

E.2. EOT and BPDA attacks

We take the Backward Pass Differentiable Approximation (BPDA) attack (Athalye et al., 2018a) for potential gradient risks and designs attacks. We design four different attacks as follows:

- BPDA₁: For potential gradients vanishing risk (i.e., small gradient of ensemble from different gradient of base learner), we instead use the k times of the average logits of the base learners as the logit of the ensemble, where $k \in [1, m]$ and m is the number of base learners. We evaluate all possible values of $k \in [1, m]$ and report the lowest adversarial accuracy observed.
- BPDA₂: For potential incorrect gradients from random gradients of single base learner, we consider the attack of deleting a base learner and using other gradients of the ensemble.
- BPDA₃: For potential incorrect gradients from random gradients of base learners, we could consider the attack of selecting randomly half of base learners at each step for attack.
- BPDA₄: For other potential gradient risks, we consider the black box attack (Andriushchenko et al., 2020).

We also take the Expectation over Transformation (EOT) attack (Athalye et al., 2018b) to add adversarial perturbations insensitively in transformations. We implement 20 times rotations randomly within -30 to +30 degrees.