

Inducing Data Amplification Using Auxiliary Datasets in Adversarial Training

Saehyung Lee* Hyungyu Lee

Department of Electric and Computer Engineering
Seoul National University, Seoul 08826, South Korea

{halo8218, rucy74}@snu.ac.kr

Abstract

Several recent studies have shown that the use of extra in-distribution data can lead to a high level of adversarial robustness. However, there is no guarantee that it will always be possible to obtain sufficient extra data for a selected dataset. In this paper, we propose a biased multi-domain adversarial training (BiaMAT) method that induces training data amplification on a primary dataset using publicly available auxiliary datasets, without requiring the class distribution match between the primary and auxiliary datasets. The proposed method can achieve increased adversarial robustness on a primary dataset by leveraging auxiliary datasets via multi-domain learning. Specifically, data amplification on both robust and non-robust features can be accomplished through the application of BiaMAT as demonstrated through a theoretical and empirical analysis. Moreover, we demonstrate that while existing methods are vulnerable to negative transfer due to the distributional discrepancy between auxiliary and primary data, the proposed method enables neural networks to flexibly leverage diverse image datasets for adversarial training by successfully handling the domain discrepancy through the application of a confidence-based selection strategy. The pre-trained models and code are available at: <https://github.com/Saehyung-Lee/BiaMAT>.

1. Introduction

The usefulness of adversarial examples in training deep neural networks (DNNs) demonstrates that the method through which these structures perceive the world is markedly different from that employed by humans. Many approaches [1] have been proposed to bridge the gap in adversarial robustness between humans and DNNs. Among these, training based on the use of adversarial examples as training data is considered as the most effective method to improve the robustness of DNNs. Unfortunately, as demon-

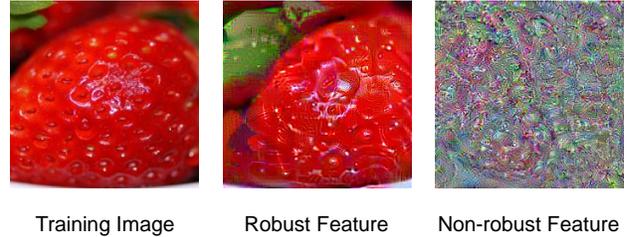


Figure 1. Visualization of robust and non-robust features [14, 21]

strated by Schmidt et al. [34], the sample complexity of adversarially robust generalization is substantially higher than that of standard generalization. To address this issue, several recent studies [4, 38] leveraged extra (in-distribution) unlabeled data and developed methods for improving the sample complexity of robust generalization. However, although such methods enable state-of-the-art adversarial robustness, they are not always capable of obtaining extra in-distribution data for any selected data distribution. In this paper, we propose a biased multi-domain adversarial training (BiaMAT) method to improve the adversarial robustness of a classifier on a primary dataset based on the use of publicly available (labeled) auxiliary datasets. The proposed method yields the desired effect based on the following assumption:

Assumption 1. *A common robust and non-robust feature space exists between the primary and auxiliary datasets.*

Figure 1 shows that robust features [21] exhibit human-perceptible patterns. We may assume that if two datasets are similar from a human perspective, then they share robust features. However, non-robust features are imperceptible to humans, thus, we cannot determine whether Assumption 1 is correct by a human perception. Fortunately, recent studies [28, 24] have provided empirical evidence in support of the presence of a common non-robust feature space among diverse image datasets. Therefore, unlike existing state-of-the-art methods [4, 38], which employ in-distribution data, under BiaMAT, the distribution of the auxiliary dataset and the corresponding primary dataset can differ. For example,

*Correspondence to: Saehyung Lee halo8218@snu.ac.kr.

by applying BiaMAT, we can leverage CIFAR-100 [22], Places365 [44], or ImageNet [8, 12] as an auxiliary dataset for adversarial training on CIFAR-10 [22].

The proposed method achieves an inductive transfer between adversarial training on the primary dataset (referred to as the “primary task”) and auxiliary dataset (referred to as “auxiliary task”). In other words, BiaMAT learns primary and auxiliary tasks in parallel within the framework of multi-domain learning [27], and the inductive bias provided by the auxiliary tasks is transferred to the primary task through a common hidden structure. This mechanism can be considered to be an increase in the size of the training dataset [5]. In addition, based on studies that have demonstrated the presence of non-robust features [40, 21], we classify the effects of adversarial training into two types and demonstrate the usefulness of the proposed method irrespective of the type considered. In particular, we dissociate the compound effect of the proposed method into the effects of *non-robust feature regularization* and *robust feature learning* and assess the contribution of each through the use of the expectation of random labels [5, 24]. Our experimental results on the CIFAR datasets and ImageNet demonstrate that BiaMAT can effectively use training signals generated from various auxiliary datasets. Furthermore, we show that while existing methods are vulnerable to negative transfer due to the distributional discrepancy between auxiliary and primary data, the proposed method enables neural networks to flexibly leverage diverse image datasets for adversarial training by successfully dealing with domain discrepancy through the application of a confidence-based selection strategy.

2. Biased multi-domain adversarial training

2.1. Method

Existing state-of-the-art methods leveraging extra unlabeled data [4, 38] first remove out-of-distribution (OOD) data from a given auxiliary dataset, and then use the remaining data for training with pseudo-labels. Therefore, they are only effective when a given auxiliary dataset contains a large number of in-distribution data, and are suboptimal in terms of data utility. To maximize the data utility, we leverage given auxiliary data via multi-domain learning [13]. Multi-domain learning is a strategy for improving the performance of tasks that solve the same problem across multiple domains by sharing information across these domains. In standard settings, domains typically share semantic features. For example, in the Office dataset [33], data that belong to the same class (*e.g.*, keyboard) are separated into different domains (*e.g.*, amazon, webcam). In adversarial settings, by contrast, it is possible to find evidence for tighter-than-expected relationships between different datasets [28]. In particular, the use of domain-agnostic adversarial examples [28] and robust training methods that leverage different datasets [6, 24] demon-

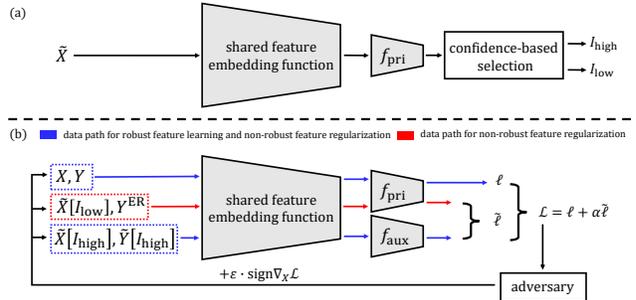


Figure 2. **Overview of BiaMAT.** We use a shared feature embedding function for both primary and auxiliary tasks, while prediction functions (f_{pri} and f_{aux}) are separated for each task. (a) At the first stage of BiaMAT, a confidence-based selection strategy classify given auxiliary data \tilde{X} into two groups: auxiliary data that share robust and non-robust features with the primary dataset ($\tilde{X}[I_{\text{high}}]$); and that share only non-robust features with the primary dataset ($\tilde{X}[I_{\text{low}}]$). (b) We adversarially train the model on $\tilde{X}[I_{\text{high}}]$ and the primary data X using their ground truth labels ($\tilde{Y}[I_{\text{high}}]$ and Y) to achieve robust feature learning and non-robust feature regularization (indicated in blue). For $\tilde{X}[I_{\text{low}}]$, however, we use the expectation of random labels y^{ER} to induce data amplification only for non-robust feature regularization (indicated in red)

strates that common adversarial spaces can exist across different datasets. Therefore, our proposed method expands the range of related domains relative to that considered under standard settings with the goal of *maximizing the adversarial robustness of the classifier on one primary dataset*. In this respect, BiaMAT differs from standard multi-domain learning, for which the primary goal is increasing the average performance over multiple domains.

We classify given auxiliary data into two types: (i) auxiliary data that share robust and non-robust features [21] with the primary dataset; and (ii) that share only non-robust features with the primary dataset. We do not consider auxiliary data sharing only robust features with the primary dataset based on studies demonstrating the presence of a common non-robust feature space among diverse image datasets [36, 24]. As previous multi-domain learning studies observed [5, 27], using (i) will be beneficial to the primary task performances. By contrast, (ii) can help the primary task by regularizing the shared non-robust features but, at the same time, it can suppress the advantages of multi-domain learning by generating an inductive bias toward extraneous robust features—an effect, called “negative transfer”. Therefore, we need a method that can acquire training signals for non-robust features without achieving inductive bias for robust features from (ii). In Sec. 2.2, we theoretically demonstrate that for each of the two data types, multi-domain learning can improve the adversarial robustness of the primary task in a simple Gaussian model. Especially, in the theoretical case of infinite batch size, we demonstrate that

the adversarial training on (ii) with random labels can improve adversarial robustness while avoiding negative transfer. Based on our theoretical analysis, we propose the use of the expectation of random labels (y^{ER}) on (ii). Interestingly, the use of y^{ER} is identical to the recently proposed OAT method [24], which uses OOD data to improve adversarial robustness. We compare our method with OAT in Appendix D. To classify given auxiliary data into the two types, our proposed method utilizes a confidence-based data selection strategy. Confidence scores are used in many research fields, including semi-supervised learning [37, 32] and OOD detection [18]. For example, in the case of semi-supervised learning, high-confidence unlabeled data are used with their pseudo-labels, while the remaining unlabeled data are not used. By contrast, in our proposed method, the use of y^{ER} naturally separates the auxiliary data according to their robust features, and adversarial training algorithms for each of the two data types are applied. We provide the detailed description of the confidence-based selection strategy in Sec. 2.4.

To sum up, we consider a multi-domain learning problem on a primary dataset $D \subset \mathcal{X} \times \mathcal{Y}$ and an auxiliary dataset $\tilde{D} \subset \mathcal{X} \times \tilde{\mathcal{Y}}$. \tilde{D} is divided into two data subsets \tilde{D}_{high} (high confidence) and \tilde{D}_{low} (low confidence) through the confidence-based selection strategy. The hypotheses of the tasks are denoted by $h_{\text{pri}} : \mathcal{X} \rightarrow \mathcal{Y}$ and $h_{\text{aux}} : \mathcal{X} \rightarrow \tilde{\mathcal{Y}}$, where $h_{\text{pri}} = f_{\text{pri}} \circ g$ and $h_{\text{aux}} = f_{\text{aux}} \circ g$. Here, f_{pri} and f_{aux} are the prediction functions that output class probabilities for the primary and auxiliary datasets, respectively, and g is the shared feature embedding function. The loss function for D is defined as $\ell = \mathbb{E}_{(\mathbf{x}, y) \sim D} [\ell_{\text{adv}}(\mathbf{x}, y; h_{\text{pri}}, S)]$, where ℓ_{adv} is the adversarial loss, and S represents the set of perturbations an adversary can apply. Any existing adversarial losses [25, 43] can be employed for ℓ_{adv} . In addition, the loss for \tilde{D} is defined as $\tilde{\ell} = \frac{1}{|\tilde{D}|} \sum_{(\tilde{\mathbf{x}}, \tilde{y}) \sim \tilde{D}_{\text{high}}} \ell_{\text{adv}}(\tilde{\mathbf{x}}, \tilde{y}; h_{\text{aux}}, S) + \frac{1}{|\tilde{D}|} \sum_{(\tilde{\mathbf{x}}, \tilde{y}) \sim \tilde{D}_{\text{low}}} \ell_{\text{adv}}(\tilde{\mathbf{x}}, y^{\text{ER}}; h_{\text{pri}}, S)$. Our goal is to attain a small adversarial loss on the primary dataset. Thus, the proposed method minimizes the following loss:

$$\mathcal{L} = \ell + \alpha \tilde{\ell}, \quad \text{where } \alpha \in [0, 1]. \quad (1)$$

α is a hyperparameter that biases the multi-domain learning toward the primary task. Although we describe BiaMAT at the dataset level, we actually apply BiaMAT at the mini-batch level. Figure 2 provides an overview of BiaMAT.

2.2. Theoretical motivation

We analyze the proposed method from the perspective of *non-robust feature regularization* and *robust feature learning*, which are the two effects of adversarial training. In particular, we define a simple Gaussian model to demonstrate how the proposed method induces training data amplification using an auxiliary dataset that satisfies Assumption 1.

Preliminary. Tsipras et al. [40] described the effect of adversarial training by constructing a classification task through which training examples $(\mathbf{x}, y) \in \mathbb{R}^{d+1} \times \{\pm 1\}$ are drawn from a distribution, as follows:

$$y \stackrel{u.a.r.}{\sim} \{-1, +1\}, \quad x_1 = \begin{cases} +y & \text{w.p. } p \\ -y & \text{w.p. } 1 - p \end{cases}, \quad (2)$$

$$x_2, \dots, x_{d+1} \stackrel{i.i.d.}{\sim} \mathcal{N}(\eta y, 1),$$

where x_1 is a robust feature that robustly correlates to the label ($p \geq 0.5$), and the remaining features x_2, \dots, x_{d+1} are non-robust features that are vulnerable to adversarial attacks ($0 < \eta < \ell_{\infty}$ -bound). For this data distribution, the authors demonstrated that the following linear classifier could attain a standard accuracy arbitrarily close to 100%, although it is susceptible to adversarial attacks:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}_{\text{unif}}^{\top} \mathbf{x}), \quad \text{where } \mathbf{w}_{\text{unif}} = \left[0, \frac{1}{d}, \dots, \frac{1}{d} \right]. \quad (3)$$

Notably, the following lemma indicates the importance of adversarial training:

Lemma 1. (Tsipras et al.) *Adversarial training results in a classifier that assigns zero weight to non-robust features x_2, \dots, x_{d+1} .*

Lemma 1 shows that adversarial training (i) lowers the sensitivity of the classifier to non-robust features and (ii) achieves a certain level of classification accuracy by learning robust features. We refer to (i) and (ii) as *non-robust feature regularization* and *robust feature learning*, respectively.

Setup and overview. Given a shared feature embedding function $g : \mathcal{X} \rightarrow \mathcal{Z}$, we define primary and auxiliary data models in the feature space \mathcal{Z} , sampled from each of the following distributions:

$$\begin{aligned} \text{(Primary)} \quad & y \stackrel{u.a.r.}{\sim} \{-1, +1\}, \quad z_1 \sim \mathcal{N}(y, u^2), \\ & z_2, \dots, z_{d+1} \stackrel{i.i.d.}{\sim} \mathcal{N}(\eta y, 1), \\ \text{(Auxiliary)} \quad & \tilde{y} = \text{sign}(\gamma) \cdot y, \quad \tilde{z}_1 \sim \mathcal{N}(y|\gamma|, v^2), \\ & \tilde{z}_2, \dots, \tilde{z}_{d+1} \stackrel{i.i.d.}{\sim} \mathcal{N}(\eta y|\gamma|, 1), \end{aligned} \quad (4)$$

where $\gamma \in [-1, 1]$ is the correlation coefficient between the two tasks. We use the accent (tilde) to represent variables associated with the ‘‘auxiliary task’’. In Eq. (4), z_1 is a robust feature that robustly correlates to the label, whereas the other features z_2, \dots, z_{d+1} are non-robust features that are vulnerable to adversarial attacks ($0 < \eta < \ell_{\infty}$ -bound). If the two datasets are highly correlated in terms of robust and non-robust features, from Eq. (3), it is evident that the following linear classifiers can achieve a high standard accuracy on the

primary and auxiliary datasets, respectively, although they have low adversarial robustness:

$$\begin{aligned} \text{(Primary)} \quad p(y | \mathbf{z}) &= \frac{1-y}{2} + y\sigma(\mathbf{w}^\top \mathbf{z}), \\ \text{(Auxiliary)} \quad p(\tilde{y} | \tilde{\mathbf{z}}) &= \frac{1-\tilde{y}}{2} + \tilde{y}\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}), \end{aligned} \quad (5)$$

where $\mathbf{w} = [0, \frac{1}{d}, \dots, \frac{1}{d}]$, and $\sigma(\cdot)$ denotes a sigmoid function. To study the effect of adversarial training on the auxiliary task on the non-robust feature regularization for the primary task, we derive the gradients of the primary and auxiliary adversarial losses with respect to the non-robust features, which are then back-propagated through the shared feature embedding function g . In addition, we demonstrate how the use of random labels enables us to dissociate the compound effect of the proposed method into the effects of non-robust feature regularization and robust feature learning.

Non-robust feature regularization. First, we generate the adversarial feature vector of our classification model ($\tilde{\mathbf{z}}^{\text{adv}} = g(\tilde{\mathbf{x}} + \boldsymbol{\delta}) : \boldsymbol{\delta} \in S$, where $\tilde{\mathbf{x}} \in \mathcal{X}$ denotes the auxiliary input vector). The objective function of the adversary to deceive our model is the cross-entropy loss [25].

Lemma 2. *Let $i \in \{2, \dots, d+1\}$ and $\eta < \lambda < 1$. Then, the expectation of the adversarial feature vector against the auxiliary task is*

$$\mathbb{E}[\tilde{z}_1^{\text{adv}}] = y, \quad \mathbb{E}[\tilde{z}_i^{\text{adv}}] = (\eta - \lambda)y. \quad (6)$$

Proof is in Appendix A. The stochastic gradient descent to the cross-entropy loss on $\tilde{\mathbf{z}}^{\text{adv}}$ is applied to update our classification model. In particular, by deriving the auxiliary loss gradient with respect to the adversarial feature vector, we determine the training signals that are generated from the auxiliary task and transferred to the primary task through the shared feature embedding function.

Theorem 1. *Let $\ell(\cdot; \mathbf{w})$ and $\tilde{\ell}(\cdot; \gamma\mathbf{w})$ be the loss functions of the primary and auxiliary tasks, respectively, and $t = \frac{1}{2}(y+1)$. When the auxiliary data are closely related to the primary data from the perspective of robust and non-robust features, i.e., $|\gamma| = 1$, the expectation of the gradient of $\tilde{\ell}$ with respect to $\tilde{z}_i^{\text{adv}} : i \in \{2, \dots, d+1\}$ is*

$$\mathbb{E}\left[\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}}\right] = \frac{1}{d}\mathbb{E}[\sigma(\mathbf{w}^\top \mathbf{z}^{\text{adv}}) - t] = \mathbb{E}\left[\frac{\partial \ell}{\partial z_i^{\text{adv}}}\right]. \quad (7)$$

The theoretical results in the cases of $|\gamma| < 1$ (weak correlation) are discussed in Appendix A. From Lemma 2 and Theorem 1, for $i \in \{2, \dots, d+1\}$, it can be seen that $\text{sign}(\mathbb{E}[\tilde{z}_i^{\text{adv}}]) = \text{sign}\left(\mathbb{E}\left[\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}}\right]\right)$. That is, the application of a gradient descent guides the shared feature embedding

function to pay less attention to non-robust features. In addition, Theorem 1 shows that if the auxiliary task is closely related to the primary task in terms of non-robust features, the training signals obtained from the auxiliary adversarial loss and back-propagated to the shared feature embedding function have the same effect as those of the primary task from the perspective of non-robust feature regularization. Therefore, this can be considered data amplification for non-robust feature regularization.

Robust feature learning. If $|\gamma| = 1$ and the weight value for the robust feature z_1 is non-zero, clearly, the auxiliary task on $\tilde{\mathbf{z}}^{\text{adv}}$ can induce data amplification for robust feature learning as well as non-robust feature regularization. However, when the auxiliary dataset contains extraneous robust features, the learning on $\tilde{\mathbf{z}}^{\text{adv}}$ may lead to negative transfer, which suppresses the advantages of multi-domain learning. To prevent inductive transfer between tasks, Caruana [5] shuffled the class labels among all samples in the auxiliary dataset. Similarly, to avoid negative transfer, the use of random labels can be considered in our case. To investigate the effect of adversarial training on the shuffled auxiliary dataset, we replace the true labels in the auxiliary data defined in Eq. (4) with random labels. That is, we define the auxiliary feature-label pairs, $(\tilde{\mathbf{z}}, q) \in \mathbb{R}^{d+1} \times \{\pm 1\}$, sampled from a distribution as follows:

$$\begin{aligned} \tilde{z}_1 &\sim \mathcal{N}(y|\gamma|, v^2), \quad \tilde{z}_2, \dots, \tilde{z}_{d+1} \stackrel{i.i.d.}{\sim} \mathcal{N}(\eta y|\gamma|, 1), \\ q &\stackrel{u.a.r.}{\sim} \{-1, +1\}. \end{aligned} \quad (8)$$

For the case in which the auxiliary task is adversarially trained on $(\tilde{\mathbf{z}}, q)$ pairs with the cross-entropy loss, the following theorem can be proven:

Theorem 2. *Let $\tilde{\ell}(\cdot; \gamma\mathbf{w})$ be the loss function of the auxiliary task. Then, if $|\gamma| = 1$, with high probability, the signs of $\tilde{z}_i^{\text{adv}} : i \in \{2, \dots, d+1\}$ and the auxiliary loss gradient with respect to \tilde{z}_i^{adv} are*

$$\text{sign}(\tilde{z}_i^{\text{adv}}) = -\gamma q = \text{sign}\left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}}\right). \quad (9)$$

Because the gradient with respect to \tilde{z}_i^{adv} is of the same sign as \tilde{z}_i^{adv} with high probability, the application of a gradient descent makes the shared feature embedding function to refrain from using non-robust features, thereby enabling the model to achieve non-robust feature regularization. Conversely, the shuffled dataset cannot provide any robust features because the use of random labels completely eliminates the relationship between images and labels. To further investigate the effect of adversarial training on the shuffled auxiliary dataset with regard to robust feature learning, we assign a positive number to the weight (defined in Eq. (5))

Table 1. Accuracy (under Autoattack [10]) comparison of the models adversarially trained [25] on CIFAR-10 using either \tilde{y} or y^{ER} for the auxiliary dataset. The best result for each auxiliary dataset is indicated in bold; auxiliary datasets that produce better results when used with y^{ER} are indicated in red

\tilde{y}/y^{ER}	Auxiliary dataset			ImageNet	Baseline [25]
	SVHN	CIFAR-100	Places365		
\tilde{y}	47.44	48.48	48.88	50.33	48.53
y^{ER}	48.53	49.89	49.24	49.81	

corresponding to the robust feature z_1 and derive the training signals that are sent to the shared feature embedding function as follows:

Theorem 3. Let $\tilde{\ell}(\cdot; \gamma w)$ be the loss function of the auxiliary task. Then, if $|\gamma| = 1$ and $w_1 > 0$, with high probability, the signs of \tilde{z}_1^{adv} and the auxiliary loss gradient with respect to \tilde{z}_1^{adv} are

$$\text{sign}(\tilde{z}_1^{\text{adv}}) = y, \quad \text{sign}\left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1^{\text{adv}}}\right) = -\gamma q. \quad (10)$$

Assuming that the classification model is still vulnerable to adversarial examples, $\left|\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1^{\text{adv}}}\right|$ is independent of q because an adversary can always yield a large loss regardless of q . Hence, in the theoretical case of an infinite batch size, the adversarial training on the shuffled auxiliary dataset will not affect the robust feature learning for the primary task because y and q are independent of each other, and q is sampled uniformly at random. In practice, however, the minibatch gradient descent is employed to train DNNs, and thus, unfavorable training signals can be generated from the auxiliary task on the shuffled dataset in terms of robust feature learning. To resolve this issue, we use the expectation of random labels instead of the one-hot random labels. Furthermore, to close the gap between the theory ($|\gamma| = 1$) and practice, we use a shared prediction function for the primary and the shuffled auxiliary data. In other words, we assign $y^{\text{ER}} = [\frac{1}{c}, \dots, \frac{1}{c}]$, where c is the number of the classes in the primary dataset, to the auxiliary data to observe only the effectiveness of non-robust feature regularization while excluding the contribution of robust feature learning.

2.3. Empirical evidence

To empirically demonstrate the arguments developed above, we conduct a test confirming the following two statements: (i) When \tilde{D} has a weak relationship with D in terms of robust features, the use of y^{ER} ($I_{\text{high}} = \emptyset$ and $I_{\text{low}} = \{0, 1, \dots, |\tilde{D}| - 1\}$ in Fig. 2) results in better adversarial robustness than that produced by the use of \tilde{y} ($I_{\text{high}} = \{0, 1, \dots, |\tilde{D}| - 1\}$ and $I_{\text{low}} = \emptyset$ in Fig. 2), and

(ii) when \tilde{D} is closely related to D in terms of robust features, the use of y^{ER} results in worse adversarial robustness than that produced by the use of \tilde{y} . Table 1 lists the results of executing the test using CIFAR-10 as the primary dataset D . Here, we use SVHN [29], CIFAR-100, and Places365 as auxiliary datasets that are weakly related to CIFAR-10 in terms of robust features, based on the previous OOD detection studies [18, 35]. In addition, we use ImageNet as an auxiliary dataset that is closely related to CIFAR-10 from the perspective of robust features [17, 7]. As shown, for SVHN, CIFAR-100, and Places365, the use of y^{ER} leads to better adversarial robustness than that induced by the use of \tilde{y} , even though image-label mappings in the auxiliary datasets are disrupted. These results demonstrate that the robust feature learning induced by using all the image-label pairs ($\tilde{x}-\tilde{y}$) in each of the SVHN, CIFAR-100, and Places365 datasets is detrimental to the primary task on CIFAR-10. By contrast, for ImageNet, the fact that blocking robust feature learning using y^{ER} leads to less performance improvement indicates that beneficial inductive transfer in terms of robust feature learning can be achieved from the auxiliary task on ImageNet. That is, ImageNet shares a large number of robust features as well as non-robust features with CIFAR-10.

2.4. A confidence-based selection strategy

Our analysis demonstrate that when using an auxiliary dataset for the primary task, the optimal algorithm to be applied varies according to the relationship between the two datasets in terms of robust features. In real world scenario, however, the auxiliary dataset may contain both favorable and unfavorable robust features for the primary task. Hence, we introduce a sample-wise selection strategy in our proposed method. The selection strategy is required to classify given auxiliary data into two groups based on the robust features, and the robust features exhibit human-perceptible patterns as shown in Fig. 1. Therefore, the objective of the selection strategy is consistent with that of existing OOD detection methods [18, 35]. Hendrycks et al. [18] proposed an OOD detection method using the confidence scores of the query data. To be specific, they trained models to give OOD samples a uniform posterior. Interestingly, the use of y^{ER} is naturally connected to their proposed method. That is, the use of y^{ER} achieves non-robust feature regularization without resulting in negative transfer while at the same time lowering the confidence scores of data irrelevant to the primary task in terms of robust features. On this basis, our selection strategy classify given auxiliary data samples based on their confidence scores. The proposed confidence-based method first (i) trains a classifier from scratch on the primary dataset; (ii) after a few epochs (warm-up), sets up a threshold using a hyperparameter $\pi \in \mathbb{R}^+$ and the mean confidence of the sampled primary data to sort out the auxiliary data samples that are likely to cause negative transfer; (iii) selects the lower-

Algorithm 1 Biased multi-domain adversarial training (BiaMAT) with the confidence-based selection strategy

Require: Primary dataset D , auxiliary dataset \tilde{D} , model parameter θ , training iterations K , warmup iterations K_w , learning rate τ , hyperparameters $\alpha \in \mathbb{R}^+$ and $\pi \in \mathbb{R}^+$

```

1: for  $k = 1$  to  $K_w$  do
2:   /* Warm-up training on  $D$  */
3:   Sample a minibatch  $B \sim D$ 
4:    $\mathcal{L} \leftarrow \mathbb{E}_{(\mathbf{x}, y) \sim B} [\ell_{\text{adv}}(\mathbf{x}, y; h_{\text{pri}}, S)]$ 
5:    $\theta \leftarrow \theta - \tau \cdot \nabla_{\theta} \mathcal{L}$ 
6: end for
7: /* Confidence threshold */
8:  $\omega \leftarrow \pi \cdot \mathbb{E}_{(\mathbf{x}, y) \sim B} [\max h_{\text{pri}}(\mathbf{x})]$ 
9: for  $k = K_w + 1$  to  $K$  do
10:  Sample a minibatch pair  $B \sim D$  and  $\tilde{B} \sim \tilde{D}$ 
11:   $\ell \leftarrow \mathbb{E}_{(\mathbf{x}, y) \sim B} [\ell_{\text{adv}}(\mathbf{x}, y; h_{\text{pri}}, S)]$ 
12:   $\ell_{\text{high}}, \ell_{\text{low}} \leftarrow 0, 0$ 
13:  /* Confidence-based selection strategy */
14:  for  $\tilde{\mathbf{x}}, \tilde{y}$  in  $\tilde{B}$  do
15:    if  $\max h_{\text{pri}}(\tilde{\mathbf{x}}) < \omega$  then
16:       $\ell_{\text{low}} += \frac{1}{|\tilde{B}|} \ell_{\text{adv}}(\tilde{\mathbf{x}}, \tilde{y}^{\text{ER}}; h_{\text{pri}}, S)$ 
17:    else
18:       $\ell_{\text{high}} += \frac{1}{|\tilde{B}|} \ell_{\text{adv}}(\tilde{\mathbf{x}}, \tilde{y}; h_{\text{aux}}, S)$ 
19:    end if
20:  end for
21:   $\mathcal{L} \leftarrow \ell + \alpha(\ell_{\text{low}} + \ell_{\text{high}})$ 
22:   $\theta \leftarrow \theta - \tau \cdot \nabla_{\theta} \mathcal{L}$ 
23: end for
24: Output: adversarially robust classifier  $h_{\text{pri}} = f_{\text{pri}} \circ g$ 

```

than-threshold auxiliary data in each training batch based on their confidence scores for the primary classes; (iv) uses the low confidence data with y^{ER} and the remaining auxiliary data with \tilde{y} . The pseudo-code is provided in Algorithm 1.

3. Experimental results and discussion

3.1. Experimental setup

Datasets. We complement our analysis with experiments conducted on the CIFAR datasets and ImageNet. ImageNet is resized [8] to dimensions of 64×64 and then randomly divided into datasets that contain 100 and 900 classes, which are termed ImgNet100 and ImgNet900, respectively. SVHN [29], Places365, and ImageNet are used as auxiliary datasets. Auxiliary data that do not fit the input size of the classifier are resized to the primary data size. For instance, when CIFAR-10 is used as the primary dataset, Places365 is down-sampled to a dimension of 32×32 , and ImageNet32x32 is leveraged.

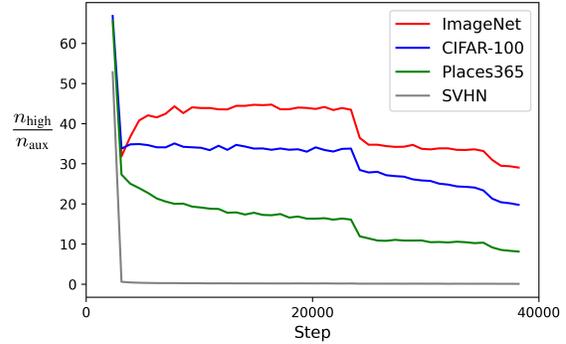


Figure 3. Ratio $\frac{n_{\text{high}}}{n_{\text{aux}}}$ for each auxiliary dataset with respect to the primary task on CIFAR-10.

Adversarial attack methods. Fast gradient sign method (FGSM) [15] is a one-step attack using the sign of the gradient. Madry et al. [25] proposed an iterative application of the FGSM method (PGD). Carlini & Wagner (CW) [3] attack is a targeted attack that maximize the logit of a target class and minimize that of ground-truth. Autoattack (AA) [10] is an ensemble attack that consists of two PGD extensions, one white-box attack [9], and one black-box attack [2]. We focus on the ℓ_{∞} -robustness, the most common robustness scenario considered in the field of heuristic defenses [25, 43, 23].

Implementation details. In our experiments, we adopt the adversarial training methods proposed by Madry et al. [25] and Zhang et al. [43] as the baseline methods, denoted by AT and TRADES, respectively. On the CIFAR datasets, we use WRN28-10 [42] and WRN34-10 for AT and TRADES, respectively. Although increasing the number of training epochs is expected to lead to higher adversarial robustness because of the use of additional data, owing to the high-computational complexity of adversarial training, we restrict the training of BiaMAT to 100 or 110 epochs with a batch size of 256 (128 primary and 128 auxiliary data samples, respectively). To evaluate adversarial robustness, we apply multiple attacks, including PGD, CW, and AA, with an ℓ_{∞} -bound with the same setting as that used in the training. PGD and CW with K iterations are denoted by PGD^K and CW^K , respectively, and the unperturbed test set is denoted by Clean. We consistently select the best checkpoint [41] to measure the adversarial robustness of the model on the test set. Further details regarding the model implementation, including an ablation study on choosing different values of π , are summarized and discussed in Appendix E.

3.2. Adversarial robustness under various attacks

Table 2 summarizes the improvements in the adversarial robustness of the models obtained from the application of BiaMAT. The proposed method can freely use various auxiliary datasets as it avoids negative transfer through the

Table 2. Performance improvements (accuracy %) on CIFAR-10, CIFAR-100, and ImgNet100 following application of the proposed method using various datasets. The best results within each baseline method (AT and TRADES) are indicated in bold

Primary dataset	Method	Auxiliary dataset	Clean	PGD ¹⁰⁰	CW ¹⁰⁰	AA
CIFAR-10	AT	-	87.37	50.87	50.93	48.53
	AT+BiaMAT (ours)	SVHN	87.34	51.90	51.40	48.61
		CIFAR-100	87.22	55.93	52.09	50.08
		Places365	87.76	57.00	51.70	49.48
		ImageNet	88.75	57.63	53.04	50.78
	TRADES	-	85.85	56.62	55.16	53.93
	TRADES+BiaMAT (ours)	SVHN	85.49	56.86	55.21	53.94
		CIFAR-100	87.02	58.69	56.85	55.48
		Places365	87.18	59.15	56.36	55.24
		ImageNet	88.03	59.80	58.01	56.64
CIFAR-100	AT	-	62.59	26.80	26.07	24.13
	AT+BiaMAT (ours)	Places365	63.44	32.61	28.53	26.49
		ImageNet	64.05	33.74	29.78	27.65
	TRADES	-	62.04	32.53	30.07	28.82
	TRADES+BiaMAT (ours)	Places365	64.58	34.38	30.72	29.24
ImageNet		65.82	36.36	33.42	31.87	
ImgNet100	AT	-	66.60	35.46	31.90	29.54
	AT+BiaMAT (ours)	Places365	70.04	40.52	33.24	30.64
		ImgNet900	68.00	40.18	35.00	32.88
	TRADES	-	56.16	27.90	22.98	21.90
TRADES+BiaMAT (ours)	Places365	57.80	29.30	24.14	23.06	
	ImageNet	58.76	31.26	25.98	24.98	

application of the confidence-based selection strategy; in fact, a comparison of Tab. 1 and Tab. 2 demonstrates that the proposed method effectively overcomes negative transfer and achieves only beneficial training signals for the primary task from the auxiliary task. To observe how the confidence-based selection strategy works while the model is being trained through the application of the proposed method, we define a ratio $\frac{n_{\text{high}}}{n_{\text{aux}}}$, where n_{aux} and n_{high} denote the amount of auxiliary data and the higher-than-threshold auxiliary data within each training batch, respectively. That is, the ratio represents the percentage of data used for robust feature learning as well as non-robust feature regularization for an auxiliary dataset. Figure 3 shows the plot of the ratio $\frac{n_{\text{high}}}{n_{\text{aux}}}$ at $\pi = 0.55$ (defined in Algorithm 1) during the training of the AT+BiaMAT models using various auxiliary datasets on CIFAR-10. As shown, the confidence-based selection strategy successfully filters out data that are likely to induce negative transfer for the primary task. In other words, a relatively high percentage of ImageNet data are used for robust feature learning, and each of the SVHN, CIFAR100, and Places365 datasets are mostly used with y^{ER} , which is con-

sistent with the results listed in Tab. 1. Additional analysis is in Appendix F.

We conduct several experiments to further investigate the proposed method. (Appendix B) To observe the effects of the use of more auxiliary datasets, we train a BiaMAT model using a combination of two auxiliary datasets; the results show that the use of more auxiliary datasets does not always lead to further improvements in adversarial robustness. In other words, the relationship (in terms of robust and non-robust features) between the primary and auxiliary datasets is more important to BiaMAT than the number of auxiliary datasets. (Appendix C) To further demonstrate that robust feature learning can be achieved from the auxiliary task in BiaMAT, we construct robust datasets [21] from the AT and AT+BiaMAT models and normally train models from scratch on each robust dataset (D^{AT} and D^{BiaMAT}); the results show that D^{BiaMAT} results in more robust models than those trained on D^{AT} , implying that BiaMAT enables DNNs to learn better robust features via inductive transfer between adversarial training on the primary and auxiliary datasets.

Table 3. Comparison (accuracy %) of the effectiveness of BiaMAT with the semi-supervised [4] and pre-training [17] methods on CIFAR-10.

Method	Auxiliary dataset	Clean	PGD ¹⁰⁰	CW ¹⁰⁰	AA
TRADES (baseline)	-	85.85	56.62	55.16	53.93
Hendrycks <i>et al.</i> [17]	CIFAR-100	80.21	45.68	44.52	42.36
	ImageNet	87.11	57.16	55.43	55.30
Carmon <i>et al.</i> [4]	CIFAR-100	82.61	54.32	51.64	50.81
	Places365	83.95	56.72	53.95	52.81
	ImageNet	85.42	57.46	54.66	53.79
	ImageNet-500k	86.02	59.49	56.43	55.63
TRADES+BiaMAT (ours)	CIFAR-100	87.02	58.69	56.85	55.48
	Places365	87.18	59.15	56.36	55.24
	ImageNet	88.03	59.80	58.01	56.64

3.3. Comparison with other related methods

Carmon et al. [4] proposed a semi-supervised learning technique where the training dataset is augmented with unlabeled in-distribution data; the main difference between this and BiaMAT is the distribution of additional data. For instance, Carmon et al. collected the in-distribution data of CIFAR-10 from 80 Million Tinyimages dataset [39] and used the unlabeled data with pseudo-labels. Therefore, no assumptions are required regarding the classes of the primary and auxiliary datasets in our scenario, but the semi-supervised method is ineffective when the primary and auxiliary datasets do not share the same class distributions. To demonstrate this, we assign pseudo-labels to the auxiliary data using a pre-trained classifier and configure each training batch (for TRADES) such that it contains the same amount of primary and pseudo-labeled data, as in [4]. In particular, we sort ImageNet based on the confidence in the CIFAR-10 classes and select the top 50k (or top 5k) samples for each class in CIFAR-10 (or CIFAR-100); this is denoted as ImageNet-500k. As shown in Tab. 3, the Carmon et al. [4] method exhibits lower effectiveness than the proposed method. Specifically, the results obtained using CIFAR-100 and Places365 demonstrate that the semi-supervised method is vulnerable to negative transfer because of the considerable domain discrepancy between the primary and auxiliary datasets.

Hendrycks et al. [17] demonstrated that ImageNet pre-training can improve adversarial robustness on the CIFAR datasets. However, the pre-training method is effective only when a dataset that has a distribution similar to that of the primary data and a sufficiently large number of samples is used. To demonstrate this, we adversarially pre-train the CIFAR-100 and ImageNet [17] models and then adversarially fine-tune them on CIFAR-10. The results in Tab. 3 demonstrate that the pre-training method is ineffective when leveraging datasets that do not satisfy the conditions mentioned above. In other words, because the effect achieved by the pre-training method arises from the reuse of features pre-

trained on a dataset that contains a large quantity of data with a distribution similar to that of the primary dataset, CIFAR-100 is not suitable for application of the CIFAR-10 task. Conversely, BiaMAT avoids such negative transfer through the application of a confidence-based strategy. That is, these results emphasize the high compatibility of the proposed method with a variety of datasets. Additional experimental results, including comparisons with OAT [24] or a generative model-based method [16], can be found in Appendix D.

4. Conclusions and future directions

In this study, we develop BiaMAT, a method that uses publicly available (labeled) auxiliary datasets to reduce the large gap between training and test errors in adversarial training. Our theoretical and empirical analysis demonstrate that the effectiveness of BiaMAT can be attributed to two factors: non-robust feature regularization and robust feature learning. In particular, we show that while existing methods are vulnerable to negative transfer due to the distributional discrepancy between auxiliary and primary data, BiaMAT can successfully overcome negative transfer through the application of a confidence-based selection strategy. In this study, however, the application of any method that can improve the performance of multi-domain learning is not considered. In addition, there is room for improvement in the effectiveness of BiaMAT with regard to the strategy used to avoid negative transfer. In future work, therefore, we will develop algorithms in which additional techniques, such as the use of adaptive weighting strategies [26], are implemented.

Acknowledgements: This work was supported by the BK21 FOUR program of the Education and Research Program for Future ICT Pioneers, Seoul National University in 2022.

References

- [1] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access*, 6:14410–14430, 2018.
- [2] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. *arXiv preprint arXiv:1912.00049*, 2019.
- [3] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.
- [4] Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, John C Duchi, and Percy S Liang. Unlabeled data improves adversarial robustness. In *Advances in Neural Information Processing Systems*, pages 11190–11201, 2019.
- [5] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [6] Alvin Chan, Yi Tay, and Yew-Soon Ong. What it thinks is important is important: Robustness transfers through input gradients. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [8] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an alternative to the cifar datasets. *arXiv preprint arXiv:1707.08819*, 2017.
- [9] Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. *arXiv preprint arXiv:1907.02044*, 2019.
- [10] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. *arXiv preprint arXiv:2003.01690*, 2020.
- [11] Luke N Darlow, Elliot J Crowley, Antreas Antoniou, and Amos J Storkey. Cinic-10 is not imagenet or cifar-10. *arXiv preprint arXiv:1810.03505*, 2018.
- [12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [13] Mark Dredze, Alex Kulesza, and Koby Crammer. Multi-domain learning by confidence-weighted parameter combination. *Machine Learning*, 79(1-2):123–149, 2010.
- [14] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Brandon Tran, and Aleksander Madry. Adversarial robustness as a prior for learned representations. *arXiv preprint arXiv:1906.00945*, 2019.
- [15] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [16] Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and Timothy A Mann. Improving robustness using generated data. *Advances in Neural Information Processing Systems*, 34, 2021.
- [17] Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using pre-training can improve model robustness and uncertainty. In *International Conference on Machine Learning*, pages 2712–2721. PMLR, 2019.
- [18] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. *arXiv preprint arXiv:1812.04606*, 2018.
- [19] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019.
- [20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [21] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. *arXiv preprint arXiv:1905.02175*, 2019.
- [22] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- [23] Saehyung Lee, Hyungyu Lee, and Sungroh Yoon. Adversarial vertex mixup: Toward better adversarially robust generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [24] Saehyung Lee, Changhwa Park, Hyungyu Lee, Jihun Yi, Jonghyun Lee, and Sungroh Yoon. Removing undesirable feature contributions using out-of-distribution data. In *International Conference on Learning Representations*, 2021.
- [25] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [26] YUREN MAO, Weiwei Liu, and Xuemin Lin. Adaptive adversarial multi-task representation learning. In *Proceedings of Machine Learning and Systems 2020*, pages 1832–1841. 2020.
- [27] Hyeonseob Nam and Bohyung Han. Learning multi-domain convolutional neural networks for visual tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4293–4302, 2016.
- [28] Muhammad Muzammal Naseer, Salman H Khan, Muhammad Haris Khan, Fahad Shahbaz Khan, and Fatih Porikli. Cross-domain transferability of adversarial perturbations. In *Advances in Neural Information Processing Systems*, pages 12885–12895, 2019.
- [29] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [30] Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning? In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 512–523. Curran Associates, Inc., 2020.
- [31] Tianyu Pang, Xiao Yang, Yinpeng Dong, Hang Su, and Jun Zhu. Bag of tricks for adversarial training. In *International Conference on Learning Representations*, 2021.
- [32] Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised

- learning. In *International Conference on Learning Representations*, 2021.
- [33] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pages 213–226. Springer, 2010.
- [34] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. In *Advances in Neural Information Processing Systems*, pages 5014–5026, 2018.
- [35] Vikash Sehwal, Mung Chiang, and Prateek Mittal. Ssd: A unified framework for self-supervised outlier detection. *arXiv preprint arXiv:2103.12051*, 2021.
- [36] Ali Shafahi, Parsa Saadatpanah, Chen Zhu, Amin Ghiasi, Christoph Studer, David Jacobs, and Tom Goldstein. Adversarially robust transfer learning. In *International Conference on Learning Representations*, 2020.
- [37] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in Neural Information Processing Systems*, 33:596–608, 2020.
- [38] Robert Stanforth, Alhussein Fawzi, Pushmeet Kohli, et al. Are labels required for improving adversarial robustness? *arXiv preprint arXiv:1905.13725*, 2019.
- [39] Antonio Torralba, Rob Fergus, and William T Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 30(11):1958–1970, 2008.
- [40] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- [41] Eric Wong, Leslie Rice, and Zico Kolter. Overfitting in adversarially robust deep learning. In *Proceedings of Machine Learning and Systems 2020*, pages 5304–5315. 2020.
- [42] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [43] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7472–7482, Long Beach, California, USA, 09–15 Jun 2019. PMLR. <https://github.com/yaodongyu/TRADES>.
- [44] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

A. Proofs

Lemma 2. Let $i \in \{2, \dots, d+1\}$ and $\eta < \lambda < 1$. Then, the expectation of the adversarial feature vector against the auxiliary task is

$$\mathbb{E} [\tilde{z}_1^{\text{adv}}] = y, \quad \mathbb{E} [\tilde{z}_i^{\text{adv}}] = (\eta - \lambda)y. \quad (8)$$

Proof. Our model comprises a non-linear feature embedding function $g : \mathcal{X} \rightarrow \mathcal{Z}$ and a linear classifier $f_{\gamma\mathbf{w}} : \mathcal{Z} \rightarrow \mathcal{Y}$. In addition, the theoretical model is based on two principles that reflect the behaviors of neural networks against adversarial examples: (i) the signs of the non-robust features $\tilde{z}_i : i \in \{2, \dots, d+1\}$ are switched by an adversary with high probability; (ii) the sign of the robust feature z_1 is not easily switched by an adversary. The objective of an adversary is to find an adversarial perturbation $\delta^* = \arg \max_{\delta \in \mathcal{S}} \tilde{\ell}(g(\tilde{x} + \delta), \tilde{y}; \gamma\mathbf{w})$. Because $f_{\gamma\mathbf{w}}$ is linear, we can easily determine the optimal adversarial direction in the feature space \mathcal{Z} using $\nabla_g \tilde{\ell}(g(\tilde{x} + \delta), \tilde{y}; \gamma\mathbf{w})$. Since the scale of the adversarial perturbation in the feature space is a problem of maximizing the convex function $\tilde{\ell}(g(\tilde{x} + \delta), \tilde{y}; \gamma\mathbf{w})$, as the scale of the perturbations increases, the situation is better from the adversarial point of view. However, these principles limit the scale range. By (i), $\lambda_i > \eta = |\mathbb{E} [\tilde{z}_i]|$, where $i \in \{2, \dots, d+1\}$; by (ii), $\lambda_1 < 1 = |\mathbb{E} [\tilde{z}_1]|$. Therefore, without loss of generality, the adversarial feature vector \tilde{z}^{adv} can be approximated by $\tilde{z} + \lambda \cdot \text{sign}(\nabla_{\tilde{z}} \tilde{\ell}(\tilde{z}, \tilde{y}; \gamma\mathbf{w}))$ (we set $\eta < \lambda = \lambda_1 = \dots = \lambda_{d+1} < 1$ for simplicity).

The loss function of the auxiliary task is formulated as

$$\begin{aligned} & \tilde{\ell}(\tilde{z}, \tilde{y}; \gamma\mathbf{w}) \\ &= -\tilde{t} \ln \sigma(\gamma\mathbf{w}^\top \tilde{z}) - (1 - \tilde{t}) \ln (1 - \sigma(\gamma\mathbf{w}^\top \tilde{z})), \end{aligned} \quad (13)$$

where $\tilde{t} = \frac{1}{2}(\tilde{y} + 1)$. Therefore,

$$\begin{aligned} \mathbb{E} [\tilde{z}_1^{\text{adv}}] &= \mathbb{E} \left[\tilde{z}_1 + \lambda \cdot \text{sign} \left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1} \right) \right] \\ &= y + \mathbb{E} [\lambda \cdot \text{sign} (\gamma w_1 (\sigma(\gamma\mathbf{w}^\top \tilde{z}) - \tilde{t}))] = y, \\ \mathbb{E} [\tilde{z}_i^{\text{adv}}] &= \mathbb{E} \left[\tilde{z}_i + \lambda \cdot \text{sign} \left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i} \right) \right] \\ &= \eta y + \mathbb{E} [\lambda \cdot \text{sign} (\gamma w_i (\sigma(\gamma\mathbf{w}^\top \tilde{z}) - \tilde{t}))]. \end{aligned} \quad (14)$$

We have

$$\begin{aligned} & \text{sign}(\gamma w_i (\sigma(\gamma\mathbf{w}^\top \tilde{z}) - \tilde{t})) \\ &= \text{sign}(w_i) \cdot \text{sign}(\sigma(\gamma\mathbf{w}^\top \tilde{z}) - \tilde{t}) = -y. \end{aligned} \quad (15)$$

Hence,

$$\mathbb{E} [\tilde{z}_i^{\text{adv}}] = \eta y - \lambda y, \quad (16)$$

where $i \in \{2, \dots, d+1\}$ and $t = \frac{1}{2}(y + 1)$. \square

Theorem 1. Let $\ell(\cdot; \mathbf{w})$ and $\tilde{\ell}(\cdot; \gamma\mathbf{w})$ be the loss functions of the primary and auxiliary tasks, respectively, and $t = \frac{1}{2}(y + 1)$. When the auxiliary data are closely related to the primary data from the perspective of robust and non-robust features, i.e., $|\gamma| = 1$, the expectation of the gradient of $\tilde{\ell}$ with respect to $\tilde{z}_i^{\text{adv}} : i \in \{2, \dots, d+1\}$ is

$$\begin{aligned} \mathbb{E} \left[\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}} \right] &= \frac{\gamma}{d} \mathbb{E} \left[\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - \gamma t - \frac{1 - \gamma}{2} \right] \\ &= \frac{1}{d} \mathbb{E} [\sigma(\mathbf{w}^\top \mathbf{z}^{\text{adv}}) - t] = \mathbb{E} \left[\frac{\partial \ell}{\partial z_i^{\text{adv}}} \right]. \end{aligned} \quad (9)$$

Proof. The expectation of the gradient of $\tilde{\ell}$ with respect to $\tilde{z}_i^{\text{adv}} : i \in \{2, \dots, d+1\}$ is

$$\mathbb{E} \left[\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}} \right] = \mathbb{E} \left[\frac{\gamma}{d} (\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - \tilde{t}) \right]. \quad (18)$$

Based on Equation 15, we obtain

$$\begin{aligned} & \mathbb{E} \left[\frac{\gamma}{d} (\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - \tilde{t}) \right] \\ &= \frac{\gamma}{d} \mathbb{E} \left[\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - \gamma t - \frac{1 - \gamma}{2} \right] \\ &= \frac{1}{d} \mathbb{E} [\sigma(\mathbf{w}^\top \mathbf{z}^{\text{adv}}) - t]. \end{aligned} \quad (19)$$

\square

Theorem 2. Let $\tilde{\ell}(\cdot; \gamma\mathbf{w})$ be the loss function of the auxiliary task. Then, if $|\gamma| = 1$, with high probability, the signs of $\tilde{z}_i^{\text{adv}} : i \in \{2, \dots, d+1\}$ and the auxiliary loss gradient with respect to \tilde{z}_i^{adv} are

$$\text{sign}(\tilde{z}_i^{\text{adv}}) = -\gamma q = \text{sign} \left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}} \right). \quad (11)$$

Proof. The gradient of $\tilde{\ell}$ with respect to $\tilde{z}_i : i \in \{2, \dots, d+1\}$ is

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i} = \frac{\gamma}{d} (\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}) - r), \quad \text{where } r = \frac{1}{2}(q + 1). \quad (20)$$

Therefore, the adversarial feature \tilde{z}_i^{adv} can be calculated as $\tilde{z}_i^{\text{adv}} = \tilde{z}_i - \lambda\gamma q$. Because $\mathbb{E} [\tilde{z}_i] = \eta y$ and $\eta < \lambda$, the sign of \tilde{z}_i^{adv} is equal to $-\gamma q$ with high probability. In addition, the gradient of $\tilde{\ell}$ with respect to \tilde{z}_i^{adv} is given as

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}} = \frac{\gamma}{d} (\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - r). \quad (21)$$

Considering the adversarial vulnerability of our classification model, we can rewrite $\sigma(\gamma\mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}})$ as $\frac{1}{2}(1 - \zeta q)$, where $\zeta \in (0, 1)$. Then,

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}} = \frac{\gamma}{d} \left(\frac{1}{2} - \frac{\zeta q}{2} - \frac{q}{2} - \frac{1}{2} \right) = \frac{-\gamma q}{2d} (1 + \zeta). \quad (22)$$

Hence, the sign of $\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}}$ is equal to $-\gamma q$ with high probability. \square

Theorem 3. Let $\tilde{\ell}(\cdot; \gamma \mathbf{w})$ be the loss function of the auxiliary task. Then, if $|\gamma| = 1$ and $w_1 > 0$, with high probability, the signs of \tilde{z}_1^{adv} and the auxiliary loss gradient with respect to \tilde{z}_1^{adv} are

$$\text{sign}(\tilde{z}_1^{\text{adv}}) = y, \quad \text{sign}\left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1^{\text{adv}}}\right) = -\gamma q. \quad (12)$$

Proof. The gradient of $\tilde{\ell}$ with respect to \tilde{z}_1 is

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1} = \gamma w_1 (\sigma(\gamma \mathbf{w}^\top \tilde{\mathbf{z}}) - r), \quad \text{where } r = \frac{1}{2}(q+1). \quad (23)$$

Assuming that the classification model is still vulnerable to adversarial examples, the adversarial feature \tilde{z}_1^{adv} is given as $\tilde{z}_1^{\text{adv}} = \tilde{z}_1 - \lambda \gamma q$. Because $\mathbb{E}[\tilde{z}_1] = y$ and $\lambda < 1$, the sign of \tilde{z}_1^{adv} is equal to y with high probability. In addition, the gradient of $\tilde{\ell}$ with respect to \tilde{z}_1^{adv} is

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1^{\text{adv}}} = \gamma w_1 (\sigma(\gamma \mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - r). \quad (24)$$

Considering the adversarial vulnerability of our classification model, $\sigma(\gamma \mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}})$ can be rewritten as $\frac{1}{2}(1 - \zeta q)$, where $\zeta \in (0, 1)$. Then,

$$\begin{aligned} \frac{\partial \tilde{\ell}}{\partial \tilde{z}_1^{\text{adv}}} &= \gamma w_1 \left(\frac{1}{2} - \frac{\zeta q}{2} - \frac{q}{2} - \frac{1}{2} \right) \\ &= \frac{-\gamma q w_1}{2} (1 + \zeta). \end{aligned} \quad (25)$$

Hence, the sign of $\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1^{\text{adv}}}$ is equal to $-\gamma q$ with high probability. \square

If we use $\mathbb{E}[q] = 0$ instead of sampled random labels q for non-robust feature regularization, the gradient of $\tilde{\ell}$ with respect to $\tilde{z}_i : i \in \{2, \dots, d+1\}$ is

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i} = \frac{\gamma}{d} \left(\sigma(\gamma \mathbf{w}^\top \tilde{\mathbf{z}}) - \frac{1}{2} \right). \quad (26)$$

Based on the high standard accuracy of our classification model, with high probability, the gradient of $\tilde{\ell}$ with respect to $\tilde{z}_i : i \in \{2, \dots, d+1\}$ can be rewritten as

$$\begin{aligned} \frac{\partial \tilde{\ell}}{\partial \tilde{z}_i} &= \frac{\gamma}{d} \left(\sigma(\gamma \mathbf{w}^\top \tilde{\mathbf{z}}) - \frac{1}{2} \right) \\ &= \frac{\gamma}{d} \left(\frac{1}{2}(1 + \zeta \gamma y) - \frac{1}{2} \right) = \frac{\zeta y}{2d}. \end{aligned} \quad (27)$$

Therefore, the adversarial feature \tilde{z}_i^{adv} can be calculated as $\tilde{z}_i^{\text{adv}} = \tilde{z}_i + \lambda y$. Because $\mathbb{E}[\tilde{z}_i] = \eta y$ and $\eta < \lambda$, the sign

of \tilde{z}_i^{adv} is equal to y with high probability. In addition, the gradient of $\tilde{\ell}$ with respect to \tilde{z}_i^{adv} is given as

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}} = \frac{\gamma}{d} \left(\sigma(\gamma \mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - \frac{1}{2} \right). \quad (28)$$

Because $\tilde{z}_i^{\text{adv}} = \tilde{z}_i + \lambda y$, $\sigma(\gamma \mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}})$ can be approximated by $\frac{1}{2}(1 + \gamma y)$. Then,

$$\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}} = \frac{y}{2d}. \quad (29)$$

Hence, with high probability, the signs of \tilde{z}_i^{adv} and the auxiliary loss gradient with respect to \tilde{z}_i^{adv} are

$$\text{sign}(\tilde{z}_i^{\text{adv}}) = y = \text{sign}\left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}}\right). \quad (30)$$

A.1. When $|\gamma| < 1$

When $|\gamma| < 1$ (weak correlation), our theorems can be replaced as follows:

Theorem 4. Let $\ell(\cdot; \mathbf{w})$ and $\tilde{\ell}(\cdot; \gamma \mathbf{w})$ be the loss functions of the primary and auxiliary tasks, respectively, and $\hat{\gamma} = \text{sign}(\gamma)$. Then, the sign of the expectation of the gradient of $\tilde{\ell}$ with respect to $\tilde{z}_i^{\text{adv}} : i \in \{2, \dots, d+1\}$ is

$$\begin{aligned} &\text{sign}\left(\mathbb{E}\left[\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}}\right]\right) \\ &= \text{sign}\left(\mathbb{E}\left[\frac{\gamma \hat{\gamma}}{d} \sigma(|\gamma| \mathbf{w}^\top \tilde{\mathbf{z}}^{\text{adv}}) - t\right]\right) = -y \\ &= \text{sign}\left(\mathbb{E}\left[\frac{1}{d} \sigma(\mathbf{w}^\top \mathbf{z}^{\text{adv}}) - t\right]\right) \\ &= \text{sign}\left(\mathbb{E}\left[\frac{\partial \ell}{\partial z_i^{\text{adv}}}\right]\right). \end{aligned} \quad (31)$$

Theorem 5. Let $\tilde{\ell}(\cdot; \gamma \mathbf{w})$ be the loss function of the auxiliary task and $\hat{\gamma} = \text{sign}(\gamma)$. Then, with high probability, the signs of $\tilde{z}_i^{\text{adv}} : i \in \{2, \dots, d+1\}$ and the auxiliary loss gradient with respect to \tilde{z}_i^{adv} are

$$\text{sign}(\tilde{z}_i^{\text{adv}}) = -\hat{\gamma} q = \text{sign}\left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_i^{\text{adv}}}\right). \quad (32)$$

Theorem 6. Let $\tilde{\ell}(\cdot; \gamma \mathbf{w})$ be the loss function of the auxiliary task and $\hat{\gamma} = \text{sign}(\gamma)$. Then, if $|\gamma| = 1$ and $w_1 > 0$, with high probability, the signs of \tilde{z}_1^{adv} and the auxiliary loss gradient with respect to \tilde{z}_1^{adv} are

$$\text{sign}(\tilde{z}_1^{\text{adv}}) = y, \quad \text{sign}\left(\frac{\partial \tilde{\ell}}{\partial \tilde{z}_1^{\text{adv}}}\right) = -\hat{\gamma} q. \quad (33)$$

The theorems in the cases of $|\gamma| < 1$ show that the scale of the correlation coefficient does not change our main idea. Moreover, the training signals generated from the auxiliary task are weakened as $|\gamma|$ approaches 0 (shown in Equation 31). Note that we consider only a common robust and non-robust feature space between the primary and auxiliary data in our theoretical model. Therefore, negative transfer, induced by learning exclusive features of auxiliary tasks, cannot be described in our model.

B. The effects of the use of more auxiliary datasets

We investigate the effects of the use of more auxiliary datasets under the proposed method and provide the experimental results in Table 4. The results demonstrate that the use of more auxiliary datasets does not always lead to further improvements in adversarial robustness. The results on CIFAR-10 indicate that the use of both SVHN and CIFAR-100 results in a lower degree of robustness than that achieved by using CIFAR-100 alone. Likewise, leveraging a combination of ImageNet and Places365 leads to more vulnerable models than that utilizing only ImageNet. In other words, the relationship between the primary and auxiliary datasets is more important to the proposed method than the number of auxiliary datasets.

In fact, this result is a general phenomenon that can be easily observed even in non-adversarial setting. To show this, we conducted an additional test in which: (1) the CIFAR-10 training set was classified into datasets that contain 25000, 12500, and 12500 samples, namely cifar-A, cifar-B, and cifar-C, respectively. We added uniform noise to the cifar-C dataset to sparsify the information included in the cifar-C dataset; (2) a classifier (ResNet18) was then trained on cifar-A using cifar-B and cifar-C as extra datasets with a batch size of 128 and evaluated on the test set. The results in Table 5 indicate that although cifar-B and cifar-C each result in performance improvement as an additional data set, the use of both cifar-B and cifar-C results in a test accuracy lower than that achieved by using cifar-B alone. We hypothesize that this is because the density of information in the training dataset is more important than the total amount of information included in the training dataset in terms of the minibatch gradient descent. In other words, when DNNs are trained with a small batch size, the quality of each minibatch gradient is more important than the total amount of information in the dataset. To confirm this, we additionally run the abovementioned experiments with larger batch sizes; in fact, Table 5 reveal that the use of both cifar-B and cifar-C results in a higher test accuracy than that achieved by using cifar-B alone in large batch settings.

C. Robust dataset analysis

Ilyas et al. [21] generated a robust dataset containing only robust features (relevant to an adversarially trained model) to demonstrate their existence in images. In particular, they optimized:

$$\min_{x_r} \|g(x_r) - g(x)\|_2$$

, where x is the target image and g is the feature embedding function. They initialized x_r as a different randomly chosen image from the training set. Thus, the robust dataset consists of optimized x_r -target label y pairs.

To confirm robust feature learning through the application of the proposed method, we construct robust datasets from the AT and AT+BiaMAT models. We then normally train models from scratch on each robust dataset using the cross-entropy loss and list the results in Table 6. As shown, the robust dataset developed using the model trained with the proposed method results in more accurate and robust models than those trained on the robust dataset of the baseline model. The proposed method thus enables neural networks to learn better robust features via inductive transfer between adversarial training on the primary and auxiliary datasets.

D. Comparison with other related methods

Semi-supervised learning. Carmon et al. [4] and Stanforth et al. [38] proposed a semi-supervised learning technique by augmenting the training dataset with unlabeled in-distribution data. The main difference between them and BiaMAT is the distribution of additional data leveraged. For instance, Carmon et al. [4] collected in-distribution data of the CIFAR-10 dataset from 80 Million Tinyimages dataset [39] and used the unlabeled data with pseudo labels. Carmon et al. [38] categorized CIFAR-10 into labeled and unlabeled data. Their theoretical analysis also assumed that the unlabeled data were in-distribution, and when out-of-distribution data were used instead, a large performance drop can be observed. Therefore, while no assumptions are required for the classes of the primary and auxiliary datasets in our scenario, the semi-supervised methods are ineffective when the primary and auxiliary datasets do not share the same class distribution. To demonstrate this, we assign pseudo labels to the auxiliary data using a classifier trained on each primary dataset and configure each training batch to contain the same amount of primary data and pseudo-labeled data as in [4]. In particular, we sort the ImageNet data based on the confidence in the primary dataset classes and select the top $(N \times 10)k$ (or top $(N \times 1)k$) samples for each class in CIFAR-10 (or CIFAR-100); this is denoted by ImageNet- $(N \times 100)k$. In Table 7, the Carmon et al. [4] method exhibits lower compatibility than the proposed method. In particular, the results obtained using CIFAR-100 and Places365 demonstrate that the semi-supervised method is vulnerable to negative trans-

Table 4. Performance improvements (accuracy %) on CIFAR-10 following application of the proposed method using various datasets. The best result is indicated in bold.

Method	Auxiliary dataset	Clean	PGD ¹⁰⁰	CW ¹⁰⁰	AA
AT	-	87.37	50.87	50.93	48.53
AT+BiaMAT	SVHN	87.34	51.90	51.40	48.61
	CIFAR-100	87.22	55.93	52.09	50.08
	SVHN, CIFAR-100	87.61	54.58	52.03	49.88
	Places365	87.76	57.00	51.70	49.48
	ImageNet	88.75	57.63	53.04	50.78
	Places365,ImageNet	87.88	56.22	51.86	49.58

Table 5. Comparison (accuracy %) of the effectiveness of data augmentation (cifar-B and cifar-C) on cifar-A.

Batch size	Dataset	Test error (mean±std over 5 runs)
128	cifar-A	9.58±0.21
	cifar-A + cifar-B	7.32±0.14
	cifar-A + cifar-C	9.15±0.26
	cifar-A + cifar-B +cifar-C	7.45±0.21
256	cifar-A	10.48±0.21
	cifar-A + cifar-B	8.06±0.18
	cifar-A + cifar-C	9.78±0.25
	cifar-A + cifar-B +cifar-C	8.12±0.20
384	cifar-A	11.08±0.35
	cifar-A + cifar-B	8.58±0.22
	cifar-A + cifar-C	10.70±0.25
	cifar-A + cifar-B +cifar-C	8.29±0.21
512	cifar-A	11.49±0.20
	cifar-A + cifar-B	9.22±0.12
	cifar-A + cifar-C	11.21±0.27
	cifar-A + cifar-B +cifar-C	8.94±0.20
1024	cifar-A	13.22±0.25
	cifar-A + cifar-B	10.55±0.21
	cifar-A + cifar-C	12.85±0.33
	cifar-A + cifar-B +cifar-C	10.23±0.17

Table 6. Accuracy (%) comparison of the models (WRN34-10) trained on each robust dataset generated from the AT and AT+BiaMAT models.

Source model	Clean	FGSM (mean±std over 5 runs)
AT	87.49±0.20	30.79±1.16
AT+BiaMAT	88.19±0.16	31.82±1.06

fer because of the considerable domain discrepancy between the primary and auxiliary datasets.

Pre-training. Hendrycks et al. [17] demonstrated that ImageNet pre-training can significantly improve adversarial

robustness on CIFAR-10. Although adversarial training on ImageNet is expensive, fine-tuning on the primary dataset does not require an extensive number of computations once the pre-trained model has been acquired. However, once this has been done, it is difficult to obtain benefit from the application of cutting-edge methods in the fine-tuning phase because the hypothesis converges in the same basin in the loss landscape [30] when trained from pre-trained weights. For example, as shown in Table 2, TRADES generally achieves higher adversarial robustness than AT. However, fine-tuning a pre-trained ImageNet model [17] through AT and TRADES, respectively, produces two models that exhibit similar levels of adversarial accuracy on CIFAR-10 (see Table 8). By contrast, the proposed method can directly bene-

Table 7. Comparison (accuracy %) of the effectiveness of BiaMAT with the semi-supervised [4] and pre-training [17] methods on the CIFAR datasets.

Primary dataset	Method	Auxiliary dataset	Clean	AA
CIFAR-10	Hendrycks et al. [17]	CIFAR-100	80.21	42.36
		ImageNet	87.11	55.30
	Carmon et al. [4]	CIFAR-100	82.61	50.81
		Places365	83.95	52.81
		ImageNet	85.42	53.79
		ImageNet-500k	86.02	55.63
		ImageNet-250k	86.51	56.27
	Gowal et al. [16]	ImageNet-100k	86.87	56.56
		Generated data [20]	85.07	57.62
	TRADES+BiaMAT (ours)	CIFAR-100	87.02	55.48
Places365		87.18	55.24	
ImageNet		88.03	56.64	
CIFAR-100	Hendrycks et al. [17]	ImageNet	59.23	28.79
		Places365	56.74	26.22
	Carmon et al. [4]	ImageNet	63.45	27.71
		ImageNet-500k	64.90	28.64
		ImageNet-250k	66.18	29.49
		ImageNet-100k	65.40	30.61
	Gowal et al. [16]	Generated data [20]	60.66	29.94
		TRADES+BiaMAT (ours)	Places365	64.58
	ImageNet		65.82	31.87

Table 8. Comparison (accuracy %) of the effectiveness of pre-training-based method using pre-trained ImageNet model on CIFAR-10 according to fine-tuning method

Fine-tuning	Clean	PGD20	PGD100
AT	87.11	57.29	56.99
TRADES	83.97	57.17	57.07

fit from the application of state-of-the-art adversarial training methods [43, 4]. BiaMAT does not require complex operations and can also leverage a variety of datasets, whereas the pre-training method is effective only when a dataset that has a distribution similar to that of the primary dataset and a sufficiently large number of samples is used. To demonstrate this difference empirically, we adversarially pre-train the CIFAR-100 and ImageNet models and then adversarially fine-tune them on CIFAR-10. The results in Table 7 demonstrate that the pre-training method is ineffective when leveraging datasets that do not satisfy the conditions mentioned above. In other words, because the effect achieved by the pre-training method arises from the reuse of features pre-trained on a dataset that contains a large quantity of data with a distribution similar to that of the primary dataset, CIFAR-100 are not suitable for application of the CIFAR-

10 task. Conversely, BiaMAT avoids such negative transfer through the application of a confidence-based selection strategy. That is, these results emphasize the high compatibility of the proposed method with a variety of datasets.

Out-of-distribution data augmented training. Out-of-distribution data augmented training (OAT) [24] was proposed as a means of supplementing the training data required for adversarial training. Under the assumption that non-robust features are shared among different datasets, the authors theoretically demonstrated that using out-of-distribution data with a uniform distribution label can reduce the contribution of non-robust features and empirically demonstrated that their method promotes the adversarial robustness of a model. OAT is similar to our proposed method in that it improves adversarial robustness by using additional data with a distribution that differs from that of the primary data. However, OAT does not derive useful information in terms of robust feature learning from auxiliary datasets. This is because OAT can only eliminate the contribution of features from the auxiliary dataset. Therefore, BiaMAT outperforms OAT when the auxiliary dataset has a close relationship with the primary dataset in terms of robust features. By contrast, if the auxiliary dataset contains a

Table 9. Results on CIFAR-10 when ImageNet-100k is auxiliary

Method	Clean	AA
OAT [24]	86.28	51.54
BiaMAT	88.23	57.01

Table 10. Performance improvements on CIFAR-10 (WRN16-8)

BiaMAT	Clean		AA		
	[16]	BiaMAT+[16]	BiaMAT	[16]	BiaMAT+[16]
84.51	82.68	83.71	51.48	52.74	53.21

large amount of useful information in terms of non-robust feature regularization rather than robust feature learning, the improvements resulting from the applications of OAT and BiaMAT can be similar.

BiaMAT has two advantages over OAT and RST: (i) OAT and RST assume that the given auxiliary dataset is out-distribution (OOD) and in-distribution (ID), respectively. Hence, if a dataset contains both OOD and ID samples, they need an additional filtering process. On contrary, BiaMAT is an end-to-end method that does not need any filtering; (ii) If the assumptions on auxiliary datasets do not hold, OAT and RST will perform badly. *E.g.*, OAT using ImageNet-100k (100k ImageNet samples closest to CIFAR-10) as an auxiliary dataset deteriorates the robustness on CIFAR-10. Tab. 9 indicates that in that case the BiaMAT model outperforms the OAT model by a large margin.

Generated data. Recently, Goyal et al. [16] leveraged generative models [20] to artificially increase the training dataset size. They showed that state-of-the-art robust accuracy can be achieved by using the increased training dataset. To be specific, they demonstrated that their proposed method yields the desired effect under the following conditions: (i) The pre-trained non-robust classifier (pseudo-label generator) must be accurate on all realistic inputs. (ii) The generative model accurately approximate the true data distribution. From these conditions, we can infer the limitations of their method. That is, the effectiveness of their method is highly dependent on the quality of the generative and classification models that are solely trained on the original training dataset; in fact, Tab. 7 demonstrates that the use of synthetic data leads to a significant robustness improvement on CIFAR-10 (+3.69%), whereas a much smaller robustness improvement on CIFAR-100 (+1.12%) than that induce by BiaMAT (+3.05%). In addition, Tab. 7 shows that while [16] significantly improves robustness against AA, it has no effect on Clean. Based on these, we investigate whether the combination of [16] and BiaMAT, which considerably increases Clean, has a synergistic effect. Tab. 10 indicates that BiaMAT can further improve [16].

Table 11. The training times of the models in our experiments.

Primary dataset	Method	Training time (h)
CIFAR	AT	34
	AT+BiaMAT (naive)	56
	AT+BiaMAT	56.5
	TRADES	52
	TRADES+BiaMAT	103
ImgNet100	AT	119
	AT+BiaMAT	196

E. Implementation details

In all our experiments, we employed commonly used data augmentation techniques such as random cropping and flipping. On the CIFAR datasets, we used WRN28-10 [42] and WRN34-10 for AT and TRADES, respectively. On ImgNet100, we used WRN16-10.

Datasets. The CIFAR-10 dataset [22] contains 50K training and 10K test images over ten classes. The CIFAR-100 dataset [22] includes 50K training and 10K test images over one hundred classes. Each image in CIFAR-10 and CIFAR-100 consists of 32×32 pixels. The ImageNet dataset [12] has 1,281,167 training and 100,000 test images over 1,000 classes. Chrabaszcz et al. [8] created downsampled versions of ImageNet. These datasets (ImageNet32x32 and ImageNet64x64) [8] contain the identical number of images and their classes as the original ImageNet dataset. The images therein are downsampled versions having pixel sizes of 32×32 and 64×64 , respectively. SVHN is obtained from a very large set of images from urban areas in various countries using Google Street View. The CIFAR datasets are labeled subsets of the 80 million tiny images dataset [39], and the 80 million tiny images dataset contains images downloaded from seven independent image search engines: Altavista, Ask, Flickr, Cydral, Google, Picsearch, and Webshots. The Places365 images are queried from several online image search engines (Google Images, Bing Images, and Flickr) using a set of WordNet synonyms. The ImageNet images are collected from online image search engines and organized by the semantic hierarchy of WordNet.

Training time. The training times of the models are summarized in Tables 11. We used a single Tesla V100 GPU with CUDA10.2 and CuDNN7.6.5. Because of the increased training dataset size (and batch size) in the proposed method, the training time was almost twice that of the baseline method. Furthermore, a comparison of AT+BiaMAT(naive) and AT+BiaMAT revealed that the proposed confidence-based selection strategy requires negligible time.

Table 1. For the experiments in Table 1, we executed 100 training epochs on CIFAR-10. The initial learning rate was set to 0.1, and the learning rate decay was applied at 60% and 90% of the total training epochs with a decay factor of 0.1. Weight decay factor and ℓ_∞ -bound were set to $2e-4$ and $\frac{8}{255}$, respectively.

Table 2. For the models associated with AT, we executed 100 training epochs (including 5 warm-up epochs) on CIFAR-10, CIFAR-100, and ImgNet100. The initial learning rate was set to 0.1, and the learning rate decay was applied at 60% and 90% of the total training epochs with a decay factor of 0.1. Weight decay factor and ℓ_∞ -bound were set to $2e-4$ and $\frac{8}{255}$, respectively. Based on a recent study [31], for the models associated with TRADES, we executed 110 training epochs (including 5 warm-up epochs) on the CIFAR datasets and ImgNet100. The initial learning rate was set to 0.1, and the learning rate decay was applied at the 100th epoch and 105th epoch with a decay factor of 0.1. Weight decay factor and ℓ_∞ -bound were set to $5e-4$ and 0.031, respectively.

The hyperparameter α and π for each model presented in Table 1 is summarized in Table 12. From Table 12, it can be observed that when the proposed method is applied with AT, it produces good results around $\alpha = 1.0$ and $\pi = 0.5$ regardless of the primary dataset used. However, when the proposed method is applied with TRADES, the optimal set of hyperparameters are dependent on the characteristics of the primary task, such as the scale of training loss and its learning difficulty. For example, the primary task on CIFAR-10 achieves a lower training loss than that on CIFAR-100, and thus, a smaller α value is required when the primary dataset is CIFAR-10 than that required when the primary dataset is CIFAR-100. In addition, when the proposed method is applied to improve the sample complexity of a high-difficulty task, the confidence-based selection strategy becomes sensitive to the hyperparameter π , because the threshold used by the strategy is determined based on the confidences of the sampled primary data. Therefore, as a future research direction, we aim to develop an algorithm that can stably detect the data samples causing negative transfer.

When CIFAR-10 is the primary dataset, we use the same adversarial loss function for the primary and auxiliary tasks under BiaMAT. However, this setting can be problematic when the TRADES+BiaMAT model is trained on CIFAR-100. TRADES uses the prediction of natural examples instead of labels to maximize the adversarial loss. In this respect, when an insufficient training time is applied to a challenging dataset, such as CIFAR-100 and ImageNet, low-quality training signals can arise owing to the inaccurate predictions. Therefore, in our experiment, the cross-entropy loss with labels is used for auxiliary tasks when the primary dataset is CIFAR-100. The application of the cross-entropy loss function allows the TRADES+BiaMAT models

to achieve a high level of adversarial robustness on CIFAR-100, as shown in Table 2.

Pre-training. In the pre-training phase, the model was adversarially trained on the auxiliary dataset according to the implementation details described in Section 3.1. The fine-tuning phase commenced from the best checkpoint of the pre-training phase. We adversarially fine-tuned the entire layers of the pre-trained model on the primary dataset. The learning rate was set according to the global step over the pre-training and fine-tuning phase. For example, if the best checkpoint was acquired at the 65th epoch in the pre-training phase, the learning rate of the fine-tuning phase commenced at 0.01 and decreased to 0.001 after 25 epochs. When SVHN and CIFAR-100 were used as the auxiliary datasets, the abovementioned type of learning rate schedule rendered better robustness than that achieved by fine-tuning the model with a fixed learning rate [17].

E.1. Ablation study on the hyperparameter π

Here, we provide the results of ablation study on π in Table 13. From the results of the AT+BiaMAT model, the effectiveness of BiaMAT is smooth near the optimal π when it is applied with AT. In the results of TRADES+BiaMAT, however, it can be seen that the effectiveness of the proposed method is relatively sensitive to π when it is applied with TRADES. We speculate that this is because of the relatively complex loss function of TRADES, which introduces another regularization hyperparameter β [43]. Therefore, in future work, we will develop advanced algorithms that adaptively control the threshold in BiaMAT for learning stability.

F. Additional analysis of the confidence-based selection strategy

Since robust features exhibit human-perceptible patterns, we conjecture that auxiliary data samples more related to the primary dataset classes can contribute more to robust feature learning. From this motivation, we design our algorithm to use the expectation of random labels for the less-related samples. In particular, we adopt an automatic confidence-based sample selection strategy, widely used in existing novelty detection literature [19]. To understand how the proposed confidence-based selection strategy works in practice, we analyze the ratio of samples having higher confidences than the confidence threshold (*i.e.*, ω in Algorithm 1). If a sample contributes more to learn robust features, it tends to have a higher confidence score than less contributed samples.

We use the AT+BiaMAT model in Table 2, trained on the CIFAR-10 dataset with the ImageNet auxiliary dataset. The model shows 88.75% clean accuracy and 50.78% robust accuracy on AA. Table 14 shows the average higher-than-threshold ratio (*i.e.*, the ratio of samples contribute to learn

Table 12. The hyperparameter α and π for each model in Table 2

Primary dataset	Method	Auxiliary dataset	α	π
CIFAR-10	AT+BiaMAT	SVHN CIFAR-100 Places365 ImageNet	1.0	0.55
	TRADES+BiaMAT	CIFAR-100 Places365 ImageNet	0.5	0.5
CIFAR100	AT+BiaMAT	Places365 ImageNet	1.0	0.5
	TRADES+BiaMAT	Places365 ImageNet	1.0	0.3
ImgNet100	AT+BiaMAT	Places365 ImgNet900	1.0	0.5

Table 13. The results of ablation study on π . Primary dataset: CIFAR-10; Auxiliary dataset: ImageNet.

Method	π	AA
AT+BiaMAT	0.45	49.85
	0.50	50.35
	0.55	50.78
	0.60	50.32
	0.65	50.35
	0.70	50.69
TRADES+BiaMAT	0.45	56.42
	0.50	56.64
	0.55	56.21
	0.60	54.70
	0.65	54.95
	0.70	54.04

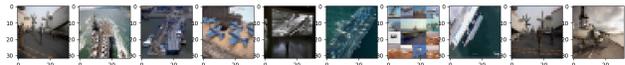
robust features) of ImageNet training images by the model. We show the average higher-than-threshold ratio for each CIFAR-10 superclasses. We match classes of two datasets by using the ImageNet synset following CINIC-10 [11]¹.

In Table 14, we observe that the related classes show higher selection ratio (larger than 50%) than the mismatched classes (29%) and the entire average (33.5%). In other words, the auxiliary samples with CIFAR-10 superclasses contribute more to robust feature learning than less related samples (“Others” in Table 14). We also illustrate the samples from the class “aircraft carrier”, showing 87.0% higher-than-threshold ratio in Figure 4. In the figure, the highest confident

¹We follow the official synset mapping used by CINIC-10 <https://github.com/BayesWatch/cinic-10/blob/master/synsets-to-cifar-10-classes.txt>



(a) The top-10 highest confident samples from “aircraft carrier” class



(b) The top-10 lowest confident samples from “aircraft carrier” class

Figure 4. The top-10 highest and lowest confident ImageNet training samples (“aircraft carrier” class) by the BiaMAT trained classifier on CIFAR-10

samples plausibly match to the CIFAR-10 superclasses, such as “Ship” and “Airplane”. On the other hand, the lowest confident samples, therefore their labels are shuffled during the training, seem to be less related to the CIFAR-10 superclasses and the original CIFAR-10 training images. The low confident samples can take a role of “out-of-distributed” dataset that can improve the confidence-based selection strategy as shown in [19].

Finally, we take a look into the “Others” classes as well. While the CIFAR-10 related classes show high higher-than-threshold ratios, we also witness that some classes not highly related to the CIFAR-10 superclasses, but weakly related to them also show high higher-than-threshold ratios. For example, (“grey whale”, 0.750), (“promontory”, 0.749), (“breakwater”, 0.734), (“dock”, 0.730), (“geyser”, 0.728), and (“sandbar”, 0.717) are not directly included in the CIFAR-10 superclasses, but share the similar environmental backgrounds (e.g., “grey whale” and “ship” are usually on the ocean background). The multi-domain learning strategy by BiaMAT let the model learn an auxiliary information by discriminating between such weakly related auxiliary classes

Table 14. Average higher-than-threshold ratio of the ImageNet training images by the AT+BiaMAT-trained CIFAR-10 classifier. The fine-grained ImageNet classes are mapped to CIFAR-10 superclasses by the WordNet hierarchy. “All” denotes the entire training ImageNet images. “Deer” and “Horse” classes has zero error because there is only one ImageNet class matched to each of them (Table ??).

CIFAR-10 Superclass	Average higher-than-threshold ratio	Standard error
Airplane	0.849	0.096
Automobile	0.706	0.163
Bird	0.554	0.143
Cat	0.501	0.136
Deer	0.720	-
Dog	0.592	0.103
Frog	0.653	0.070
Horse	0.819	-
Ship	0.677	0.215
Truck	0.763	0.129
Others (dismatched)	0.290	0.196
All	0.335	0.219

and the CIFAR-10 superclasses. Our BiaMAT can learn better robust features by the additional tasks to discriminate weak auxiliary classes from the target classes.

To sum up, our confidence-based selection strategy let the model learn better robust features from plausible extra images, while less plausible images improve the performance of the confidence-based selection strategy. At the same time, the multi-domain learning strategy by BiaMAT makes the model learn discriminative features between the samples highly correlated with target classes and the sample weakly correlated with targets (e.g., “grey whale”), thus BiaMAT shows a good robust feature learning capability. Therefore, BiaMAT can learn diverse and fine-grained features using extra images related to the target classes without suffering from the negative transfer, resulting in showing better robustness generalizability.

From these observations, we conclude that by learning robust features from extra images but related to the primary dataset, a model can learn more diverse and fine-grained features, resulting in better robustness generalizability.