# A Unified Analysis of Mixed Sample Data Augmentation: A Loss Function Perspective

Chanwoo Park \* MIT cpark97@mit.edu Sangdoo Yun \* NAVER AI Lab sangdoo.yun@navercorp.com Sanghyuk Chun NAVER AI Lab sanghyuk.c@navercorp.com

#### Abstract

We propose the first unified theoretical analysis of mixed sample data augmentation (MSDA), such as Mixup and CutMix. Our theoretical results show that regardless of the choice of the mixing strategy, MSDA behaves as a pixel-level regularization of the underlying training loss and a regularization of the first layer parameters. Similarly, our theoretical results support that the MSDA training strategy can improve adversarial robustness and generalization compared to the vanilla training strategy. Using the theoretical results, we provide a high-level understanding of how different design choices of MSDA work differently. For example, we show that the most popular MSDA methods, Mixup and CutMix, behave differently, e.g., CutMix regularizes the input gradients by pixel distances, while Mixup regularizes the input gradients regardless of pixel distances. Our theoretical results also show that the optimal MSDA strategy depends on tasks, datasets, or model parameters. From these observations, we propose generalized MSDAs, a Hybrid version of Mixup and CutMix (HMix) and Gaussian Mixup (GMix), simple extensions of Mixup and CutMix. Our implementation can leverage the advantages of Mixup and CutMix, while our implementation is very efficient, and the computation cost is almost neglectable as Mixup and CutMix. Our empirical study shows that our HMix and GMix outperform the previous state-of-the-art MSDA methods in CIFAR-100 and ImageNet classification tasks. Source code is available at https://github.com/naver-ai/hmix-gmix.

# 1 Introduction

As deep neural networks (DNNs) are data-hungry, the scale of datasets has become a foundation of modern DNN training; recent ground-breaking deep models are built upon gigantic datasets, such as 410B language tokens [6], 3.5B images [47], and 1.8B image-text pairs [33]. While such tremendously large-scale datasets are not always collectible, amplifying the dataset scale by synthesizing more data points by data augmentation techniques is common. Especially, *mixed sample data augmentation* (MSDA) [3, 11, 14, 16, 18, 23, 25, 30, 32, 32, 38, 39, 44, 45, 54, 57, 59, 61, 62, 64, 65, 67, 70, 72, 74] has become a standard technique to train a strong deep model by synthesizing a mixed sample from multiple (usually two) samples by combining both of their sample values and labels in a linear combination [74] or a cut-and-paste manner [70]. This simple idea, however, shows surprising performance enhancements in various applications, including image object recognition [16, 25, 39, 63, 70], semi-supervised learning [4, 58] self-supervised learning [35, 40, 42], noisy label training [43], meta-learning [69], semantic segmentation [10, 20], natural language understanding [24, 34], and audio processing [37, 48, 49]. Another advantage of MSDA beyond the performance improvements is that MSDA usually does not need domain-specific knowledge, such as strong image-specific [15] or audio-specific [53] transformations; hence MSDA can be universally employed

<sup>\*</sup>Equal contribution

by various applications. However, although MSDA shows excellent benefits in practice, there is still yet not enough understanding of how MSDA works well universally; *can MSDA always show better generalization and robustness than the standard training with a theoretical guarantee*? Furthermore, the design choice of MSDA can be significantly varying and the optimal design choice is still ambiguous. For example, Lee *et al.* [42] showed that in self-supervised learning, Mixup is more effective than CutMix, while Ren *et al.* [56] observed opposite results. The ambiguity is originated from the fact that we do not have a unified theoretical lens of understanding how different design choices affect the actual learning process; in short, *how are Mixup and CutMix different*?

There have been several attempts to theoretically understand Mixup, a special case of MSDA [8, 12, 75, 76]. They delve into the effect of Mixup in a loss function perspective, e.g., Mixup behaves as a regularization of the standard training [75], or in a learning theory perspective, e.g., Mixup training can provide an upper bound for the true loss [12, 69]. However, their analyses are limited to Mixup, while other MSDAs, such as CutMix, are poorly understandable through the lens of their analyses. In this paper, we extend the theoretical results of Zhang et al. [75] and Chidambaram et al. [12] to a general MSDA to provide a first unified theoretical lens for understanding how general MSDAs work by different choices of mixing strategies. We show that MSDA behaves as an input gradient and Hessian regularization (Theorem 1) as well as a regularizer for the first layer parameters; MSDA improves adversarial robustness (Theorem 3) and generalization (Theorem 4). Our theoretical results show that popular MSDA methods, such as Mixup and CutMix, behave differently in terms of regularization effects. Briefly, CutMix gives a strong regularization in the product of nearby distance pixel-level partial gradient and nearby distance Hessian of the estimated function f, while CutMix gives a weak regularization in the product of long-distance pixel-level partial gradient and long-distance Hessian of the estimated function f. In contrast, Mixup gives a regularization in gradient or Hessian of the estimated function f regardless of the pixel-level distance.

From our unified theoretical lens for MSDA, we can conclude that *there is no one-fit-all optimal MSDA fit to every data or model parameter*. In other words, the optimal mixing strategy depends on applications, datasets, and model architectures. It supports previous empirical observations that combining different MSDA methods (*e.g.*, alternatively using Mixup and CutMix during training) can outperform using only one MSDA [5, 50, 63, 71]. From these observations, we propose two simple MSDA methods that naturally generalize Mixup and CutMix, so that it can take advantage of both methods. Our first proposed method, Hybrid version of Mixup and CutMix (**HMix**), mixes two samples in both Mixup and CutMix manners; it first cut-and-paste two samples as CutMix, and then it linearly interpolates the out-of-box values of two samples as Mixup. We let HMix be able to behave as both Mixup and CutMix by introducing a stochastic control parameter. Our second proposed method, Gaussian Mixup (**GMix**), also mixes two samples in both Mixup and CutMix by introducing a stochastic control parameter. Our second proposed method, Gaussian Mixup (**GMix**), also mixes two samples in both Mixup and CutMix by introducing a stochastic control parameter. Our second proposed method, Gaussian Mixup (**GMix**), also mixes two samples in both Mixup and Cutmix manners; firstly we select a point, and then we mix two samples gradually using the Gaussian function. Our empirical results on CIFAR-100 and ImageNet show that HMix and GMix outperform the state-of-the-art MSDA methods, including Mixup, CutMix, and Stochastic Mixup & CutMix.

#### 2 A General Framework for Mixed Sample Data Augmentation (MSDA)

In this section, we define the formal definition of MSDA and notations. We define a training dataset as  $D = \{z_i = (x_i, y_i)\}_{i=1}^m$ , randomly sampled from a distribution  $\mathcal{P}_z$ . Here, z = (x, y) is the input (*e.g.*, an image) and output (*e.g.*, the target class label) pair. Then, for randomly selected two data samples,  $z_i$  and  $z_j$ , an augmented sample by MSDA,  $\tilde{z}_{i,j}^{(MSDA)}$ , is synthesized as follows

$$\begin{split} \tilde{z}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda) &= (\tilde{x}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda), \tilde{y}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda)) \\ \text{where,} \quad \tilde{x}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda) &= M(\lambda) \odot x_i + (1-M(\lambda)) \odot x_j \quad \text{and} \\ \quad \tilde{y}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda) &= N(\lambda) \odot y_i + (1-N(\lambda)) \odot y_j, \end{split}$$
(1)

where  $\lambda$  is the ratio parameter between samples, drawn from  $\mathcal{D}_{\lambda}$  (usually Beta distribution).  $\odot$  means a component-wise multiplication in vector (or matrix).  $M(\lambda)$  is a random variable conditioned on  $\lambda$  that indicates how we mix the input (*e.g.*, by linear interpolation [74] or by a pixel mask [70]).  $N(\lambda)$  denotes a random variable conditioned on  $\lambda$  that demonstrates how we combine the output. We assume that the output y can be one-dimensional data or a matrix; the former means regression or classification task, and the latter means semantic segmentation task. For the sake of simplicity, we let y be one-dimensional data:  $N(\lambda) = \lambda$  and  $\mathbb{E}[M(\lambda)] = \lambda \vec{1}$ . **Remark 1.** If the meaning is not ambiguous, then we sometimes omit  $\lambda$  (*i.e.*,  $M(\lambda)$  to M). For the sake of simplicity, we consider mixing only two samples (*i.e.*,  $\tilde{z}_{i,j}^{(MSDA)}(\lambda, 1 - \lambda)$ ), but we can similarly extend these analyses to mixing *n*-samples data augmentation [23, 32, 59]. If we combine *n*-samples, the ratio parameter will be a vector in general (See Appendix B).

**Remark 2.** As recent studies [38, 39, 67] have shown,  $M(\lambda)$  or  $N(\lambda)$  can depend on  $(z_i, z_j)$ , *e.g.*, by using a saliency map [38, 39] or the class activation map [67]. Since the proof techniques for our theoretical analysis are invariant to the choice of  $M(\lambda)$  and  $N(\lambda)$ , our proof techniques also can be applied to the dynamic MSDA methods. For simplicity, we assume that M is a random variable only depending on  $\lambda$ . In other words, we assume W as a random sample space, and  $M : W \times \Lambda \to \mathbb{R}^n$  is a measurable function. We left the theoretical analysis of dynamic methods to the future.

Now, we re-write the two most popular MSDA methods, Mixup [74] and CutMix [70], for *i*-th and *j*-th samples with  $\lambda$ , drawn from  $\mathcal{D}_{\lambda}$ , by using the proposed framework (Equation (1)) as follows:

$$\begin{aligned} \mathbf{Mixup:} \quad \tilde{z}_{i,j}^{(\mathrm{mixup})}(\lambda, 1-\lambda) &= (\tilde{x}_{i,j}^{(\mathrm{mixup})}(\lambda, 1-\lambda), \tilde{y}_{i,j}^{(\mathrm{mixup})}(\lambda, 1-\lambda)) \\ \text{where} \quad \tilde{x}_{i,j}^{(\mathrm{mixup})}(\lambda, 1-\lambda) &= \lambda x_i + (1-\lambda) x_j \quad \text{and} \\ \tilde{y}_{i,j}^{(\mathrm{mixup})}(\lambda, 1-\lambda) &= \lambda y_i + (1-\lambda) y_j. \end{aligned}$$

$$(2)$$

**CutMix:** 
$$\tilde{z}_{i,j}^{(\text{cutmix})}(\lambda, 1-\lambda) = (\tilde{x}_{i,j}^{(\text{cutmix})}(M, 1-M), \tilde{y}_{i,j}^{(\text{cutmix})}(\lambda, 1-\lambda))$$
  
where  $\tilde{x}_{i,j}^{(\text{cutmix})}(\tilde{M}^{(\text{cutmix})}, 1-\tilde{M}^{(\text{cutmix})}) = \tilde{M}^{(\text{cutmix})} \odot x_i + (1-\tilde{M}^{(\text{cutmix})}) \odot x_j$   
and  $\tilde{y}_{i,j}^{(\text{cutmix})}(\lambda, 1-\lambda) = \lambda y_i + (1-\lambda)y_j.$ 
(3)

Note that Equation (2) is equivalent to Equation (1) by putting  $M(\lambda) = \lambda \vec{1}$ . In Equation (3),  $\tilde{M}^{(\text{cutmix})}$  is a binary mask that indicates the location of the cropped box region with a relative area  $\lambda$ . Similarly, other MSDA variants can be easily formed as Equation (1) by introducing new  $M(\lambda)$  and  $N(\lambda)$ .

**Notations.** We define the loss function as  $l(\theta, z)$ , where  $\theta \in \Theta \subseteq \mathbb{R}^d$ . We define  $L(\theta) = \mathbb{E}_{z \sim \mathcal{P}_z} l(\theta, z)$  as the non-augmented population loss and  $L_m(\theta) = \frac{1}{m} \sum_{i=1}^m l(\theta, z_i)$  as the empirical loss for the non-augmented population. For a general MSDA, we can define MSDA loss as

$$L_m^{\text{MSDA}}(\theta) = \mathbb{E}_{i,j \sim \text{Unif}([m])} \mathbb{E}_{\lambda \sim \mathcal{D}_\lambda} \mathbb{E}_M l(\theta, \tilde{z}_{i,j}^{(\text{MSDA})}(\lambda, 1 - \lambda)).$$
(4)

Therefore, the Mixup and CutMix losses can be written as

$$L_m^{\text{mixup}}(\theta) = \frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E}_{\lambda \sim \mathcal{D}_\lambda} l(\theta, \tilde{z}_{i,j}^{(\text{mixup})}(\lambda, 1 - \lambda))$$
(5)  
$$L_m^{\text{cutmix}}(\theta) = \frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E}_{\lambda \sim \mathcal{D}_\lambda} \mathbb{E}_M l(\theta, \tilde{z}_{i,j}^{(\text{cutmix})}(\lambda, 1 - \lambda)),$$

where  $[m] = \{1, 2, ..., m\}$  and  $\mathcal{D}_{\lambda}$  is a distribution supported on [0, 1] with a conjugate prior. Throughout this paper, we consider  $\mathcal{D}_{\lambda}$  as  $\text{Beta}(\alpha, \beta)$ , a common selection for  $\lambda$  in practice. We define  $\mathcal{D}_X$  as the empirical distribution of the training dataset.

### **3** A Unified Theoretical Understanding of MSDA

In this section, we provide a unified theoretical lens of how MSDA works. Specifically, we follow the theoretical results for Mixup provided by Zhang *et al.* [75], where Zhang *et al.* have shown that Mixup is equivalent to the summation of the original loss function and a Mixup-originated regularization term. We will give a general approximation form for MSDA using  $\lambda \sim \text{Beta}(\alpha, \beta)$ . We also show that our analysis can be extended to *n*-sample mixed augmentation (See Appendix B) From an MSDA loss to an input gradient and Hessian regularization. We first consider the following class of loss functions for a twice differentiable prediction function  $f_{\theta}(x)$  (e.g., a softmax output of a neural network), a twice differentiable function h, and target y:

$$\mathcal{L} = \{ l(\theta, z) \, | \, l(\theta, z) = h(f_{\theta}(x)) - yf_{\theta}(x) \text{ for a twice differentiable function } h \}.$$

This function class  $\mathcal{L}$  includes the loss function induced by Generalized Linear Models (GLMs) and cross-entropy. Now, we introduce our first theoretical result that induces an MSDA loss (*i.e.*, Equation (4)) can be re-written as the summation of the original loss (the empirical loss for the non-augmented population loss,  $L_m(\theta)$ ) and input gradient-related regularization terms as follows.

**Theorem 1.** Consider a loss function  $l \in \mathcal{L}$ . We define  $\tilde{D}_{\lambda}$  as  $\frac{\alpha}{\alpha+\beta}Beta(\alpha+1,\beta) + \frac{\beta}{\alpha+\beta}Beta(\beta+1,\alpha)$ . Assume that  $\mathbb{E}_{r_x \sim \mathcal{D}_X}[r_x] = 0$ . Then, we can re-write the general MSDA loss (4) as

$$L_m^{MSDA}(\theta) = L_m(\theta) + \sum_{i=1}^3 \mathcal{R}_i^{(MSDA)}(\theta) + \mathbb{E}_{\lambda \sim \tilde{\mathcal{D}}(\lambda)} \mathbb{E}_M[(1-M)^{\mathsf{T}}\varphi(1-M)(1-M)], \quad (6)$$

where  $\lim_{a\to 0} \varphi(a) = 0$ ,

$$\mathcal{R}_{1}^{(MSDA)}(\theta) = \frac{1}{m} \sum_{i=1}^{m} (y_{i} - h'(f_{\theta}(x_{i}))) \left(\nabla f_{\theta}(x_{i})^{\mathsf{T}} x_{i}\right) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} (1 - \lambda),$$

$$\mathcal{R}_{2}^{(MSDA)}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbf{G}(\mathcal{D}_{X}, x_{i}, f, M),$$

$$\mathcal{R}_{3}^{(MSDA)}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbf{H}(\mathcal{D}_{X}, x_{i}, f, M),$$
(7)

and

$$\mathbf{G}(\mathcal{D}_{X},x_{i},f,M) = \mathbb{E}_{M}(1-M)^{\mathsf{T}}\mathbb{E}_{r_{x}\sim\mathcal{D}_{X}}\left(\nabla f(x_{i})\odot(r_{x}-x_{i})\left(\nabla f(x_{i})\odot(r_{x}-x_{i})\right)^{\mathsf{T}}\right)(1-M)$$

$$= \sum_{j,k\in coord} a_{jk}\partial_{j}f_{\theta}(x_{i})\partial_{k}f_{\theta}(x_{i})\left(\mathbb{E}_{r_{x}\sim\mathcal{D}_{X}}[r_{xj}r_{xk}]+x_{ij}x_{ik}\right),$$

$$\mathbf{H}(\mathcal{D}_{X},x_{i},f,M) = \mathbb{E}_{r_{x}\sim\mathcal{D}_{X}}\mathbb{E}_{M}(1-M)^{\mathsf{T}}\left(\nabla^{2}f_{\theta}(x_{i})\odot\left((r_{x}-x_{i})(r_{x}-x_{i})^{\mathsf{T}}\right)\right)(1-M)$$

$$= \sum_{j,k\in coord} a_{jk}\left(\mathbb{E}_{r_{x}\sim\mathcal{D}_{X}}[r_{xj}r_{xk}\partial_{jk}^{2}f_{\theta}(x_{i})]+x_{ij}x_{ik}\partial_{jk}^{2}f_{\theta}(x_{i})\right),$$
(8)

where

$$a_{jk} := \mathbb{E}_M[(1 - M_j)(1 - M_k)].$$
(9)

Proof outline of Theorem 1. Using the definition of  $\tilde{z}_{ij}$  and using the fact that the Binomial distribution and Beta distribution are in the conjugate, we can reformulate  $L_m^{(MSDA)}$ . In the process of reformulating  $L_m^{(MSDA)}$ , we should define  $\tilde{\mathcal{D}}_{\lambda}$ . Then, we can make a quadratic Taylor approximation of the loss term. Here,  $\mathbb{E}_{r_x}[r_x] = 0$  is used for not only the simplicity of the results, but also for the fact that using normalization in the dataset. Details can be found in Appendix A. We also show that Theorem 1 can be extended to *n*-sample MSDA methods (Appendix B). In this case, the combinatorial terms in quadratic multivariate Taylor approximation also come out.

How is our approximation accurate? We call  $\tilde{L}_m^{MSDA}(\theta) := L_m(\theta) + \sum_{i=1}^3 \mathcal{R}_i^{(MSDA)}$  as the approximate MSDA loss. Here, we empirically demonstrate that our quadratic approximation is almost accurate by following numerical validations in [8, 66, 75]. Specifically, we train logistic regression models on two-moons dataset [7] in two ways: (1) by using the original MSDA loss function (2) by using our approximated loss function. We employ two MSDA examples as below.

- The original Mixup *i.e.*,  $\lambda \sim \text{Beta}(1,1)$  and  $M = \lambda \vec{1}$
- Variants of CutMix *i.e.*,  $\lambda \sim \text{Beta}(1,1)$  and  $M = (m_1, m_2)$  such that  $m_i \sim \text{Bernoulli}(\lambda)$ .

Figure 1 displays the approximate loss function and the original loss function. According to empirical findings, we can conclude that the original MSDA loss is fairly close to the approximate MSDA loss.



Figure 1: Comparison of the original MSDA loss with the approximate MSDA loss function.



Figure 2: The negative M value results. the image is  $\lambda * \log + (1 - \lambda) * \operatorname{cat}$  where  $-0.75 \le \lambda \le 0.75$ 

What makes the difference between various MSDA methods? In the theorem, as we define  $E_M(1 - M) = 1 - \lambda$ ,  $\mathcal{R}_1^{(\text{MSDA})}$  is the same for every MSDA method. Namely, the difference between MSDA methods originated from  $\mathcal{R}_2^{(\text{MSDA})}$  and  $\mathcal{R}_3^{(\text{MSDA})}$ . Note that if we set  $M = \lambda \vec{1}$ , Theorem 1 indicates a Mixup loss (5), and the result is consistent with Zhang *et al.* [75]. In Equation (7) and Equation (8), we observe that  $\mathcal{R}_2$  is related to the input gradient  $\nabla f_{\theta}(x_i)$  and  $\mathcal{R}_3$  is related to input Hessian  $\nabla^2 f_{\theta}(x_i)$  with mask-dependent coefficients  $a_{jk}$  (9). In other words, different design choice of MSDA (*e.g.*, how to design M) will lead to different magnitudes of regularization on the input gradients and Hessians. Because the values of input gradients and Hessians are varying by datasets, tasks and model architecture choices, we can conclude that the optimal choice of M is dependent on the applications. We also describe how the other MSDA methods (*e.g.*, dynamic MSDAs [38, 39, 64]) can be interpreted through the lens of our unified analysis in Appendix I.

In addition, to show that MSDA behaves as a regularization on input gradients and Hessians for any desired  $a_{jk}$ , we also show that there always exists a MSDA design choice M for any desired regularization coefficient matrix  $A(\lambda) := (a_{ik}(\lambda))$  with the regular conditions.

**Theorem 2.** For the given  $\lambda$ , we assume  $A(\lambda) - (1 - \lambda)^2 \vec{1} \vec{1}^{\intercal}$  is a nonnegative definite matrix. Then we can construct a real-valued mask M that  $\mathbb{E}(1 - M_i)(1 - M_k) = a_{ik}$  for all j, k.

*Proof.* Setting  $M = 1 - \lambda + (A(\lambda) - (1 - \lambda)^2 \vec{1} \vec{1}^{\mathsf{T}})^{1/2} Z$  where Z is normal distribution, the theorem holds.

Note that, in the proof, M values are not bounded where typically we choose  $0 \le M_i \le 1$ . In other words, the theorem holds if we allow mask values out of [0, 1]. To investigate the potentiality of unbounded mask, we explore Mixup with unbounded masks in Figure 2. Although, allowing negative values to M can be beneficial, we leave a new mask design with unbounded values as a future work.

Unfortunately, as the target loss function (6) is mingled with the choice of mask M, data sample  $x_i$ , and pixel-level function gradient, the optimal choice of mixing strategy M is not achievable in the closed-form solution. Instead, Theorem 1 implies that there is no absolute superiority between the design choice of MSDA, but it depends on datasets and the target tasks, as our empirical observation is consistent with the theoretical interpretation. In Section 4, we will provide more examples of how different M affects the actual coefficients  $a_{jk}$  and the input gradients for better understanding.

Using the regularization term  $\mathcal{R}_2^{(MSDA)}$  (7), we can also provide a theoretical connection between MSDA methods and the notion of flatness where a more flat solution leads to better generalization in applications [9, 19, 21, 31, 36]. Inspired by Ma *et al.* [46], we split the parameters by  $\theta = (\theta_1, \theta_2)$ ,

and then the neural network can be represented by the form  $f_{\theta}(x) = f_{\theta_2}(\theta_1 x)$ . Therefore, we have

$$\nabla_{\theta_1} \tilde{f}_{\theta_2}(\theta_1 x) = \frac{\partial f}{\partial(\theta_1 x)} x^{\mathsf{T}}, \qquad \nabla_x \tilde{f}_{\theta_2}(\theta_1 x) = \theta_1^{\mathsf{T}} \frac{\partial f}{\partial(\theta_1 x)},$$

where  $\frac{\partial f}{\partial(\theta_1 x)}$  is the partial derivative of the first layer. Now, we have

$$((1-M) \odot \nabla f(x))^{\mathsf{T}} x = \operatorname{tr}(x((1-M) \odot \nabla f(x))^{\mathsf{T}}) = \operatorname{tr}\left(x\left(\theta_{1}^{\mathsf{T}}\frac{\partial f}{\partial \theta_{1}x} \odot (1-M)\right)^{\mathsf{T}}\right)$$
$$= \operatorname{tr}\left(x\left(\left(\frac{\partial f}{\partial \theta_{1}x}\right)^{\mathsf{T}}\theta_{1}\operatorname{diag}(1-M)\right)\right)$$
$$= \operatorname{tr}\left(\left(\nabla_{\theta_{1}}\tilde{f}_{\theta_{2}}(\theta_{1}x)\right)^{\mathsf{T}}\theta_{1}\operatorname{diag}(1-M)\right).$$

Note that the terms in G(8) can be re-written as follows

$$\sum_{j,k\in\text{coord}} \mathbb{E}_M[(1-M_j)(1-M_k)]\partial_j f_\theta(x_i)\partial_k f_\theta(x_i)(x_{ij}x_{ik}) = \mathbb{E}\left[(((1-M)\odot\nabla f(x))^{\mathsf{T}}x)^2\right].$$

In other words, by minimizing the regularization term  $\mathcal{R}_2^{(MSDA)}$ ,  $\int (\theta_1^{\mathsf{T}} \nabla_{\theta_1} \tilde{f}_{\theta_2})^2$ , *i.e.*, the regularization effect of flatness at the interpolation solution can be minimized in a sample-wise weighted manner. Therefore, the regularization term  $\mathcal{R}_2^{(MSDA)}$  also can be interpreted as a regularization of the first layer parameters and their partial derivative of f.

**Robustness and generalization properties of MSDA.** As a number of studies [46, 51, 52] have shown that regularizing input gradient and Hessian will give better robustness and generalization to the target network  $\theta$ , it can be shown that MSDA also has adversarial robustness properties and generalization properties based on Theorem 1. The full statement of Theorem 3 and Theorem 4 can be found in Appendix C and Appendix D, respectively.

**Theorem 3** (Informal). With the logistic loss function under the ReLU network, the approximate loss function of MSDA is greater than the adversarial loss with the  $\ell_2$  attack of size  $\epsilon \sqrt{d}$ .

*Proof outline of Theorem* 3. Defining adversarial loss function and using second order taylor expansion, we can prove that adversarial loss is less than MSDA loss.

**Theorem 4** (Informal). Under the GLM model and the regular conditions, and if we use MSDA in training, we have

$$L(\theta) \leq \tilde{L}_m^{(MSDA)}(\theta) + \sqrt{\frac{\mathcal{O}\left(\log(1/\delta)\right)}{n}}$$

with probability at least  $1 - \delta$ . This also holds for the MSE loss and a feature-level MSDA.

*Proof outline of Theorem 4.* MSDA regularization can be altered to the original empirical risk minimization problem with a constrained function set, and calculating Radamacher complexity of this function set gives the theorem.  $\Box$ 

In addition to Theorem 3 and Theorem 4, we can prove that the optimal solution of (4) can achieve a perfect classifier (*i.e.*, classifies every augmented sample x correctly) in the logistic classification setting by following Chidambaram *et al.* [12]. The full statement are in Appendix E.

**Summary.** Our unified theoretical lens for MSDA shows that for any MSDA method formed as Equation (4), the method satisfies that (1) it behaves as a regularizer of input gradients, Hessian, and the first layer parameters (Theorem 1); (2) there exists a mask M for any desired regularization coefficients  $a_{jk}$  (Theorem 2) (3) it achieves better adversarial robustness (Theorem 3) and generalization (Theorem 4) than the vanilla training. Interestingly, Theorem 1 shows the difference between different MSDA design choices (*e.g.*, different M, such as linear interpolation [74], cropped box [70]) will lead to different magnitudes of the input gradient regularization (7).



Figure 3: Examples generated by different MSDAs. From left to right, two original images to be mixed, Mixup, CutMix sample, HMix, and GMix. The first and the second rows show generated samples and their mixing masks M, respectively. We set  $\lambda = 0.65$  for all images and r = 0.5 for HMix.

# 4 Comparison of Different MSDA Design Choices: The Role of Masks

As we observed in the previous section, different design choices for MSDA (*i.e.*, the choice of M) affect to the degree of the regularization in Theorem 1 (*i.e.*,  $\mathcal{R}_2^{(\text{MSDA})}$  and  $\mathcal{R}_3^{(\text{MSDA})}$ ) by the relationships of pixels. In this section, we show how different MSDA methods lead to different regularization effects by empirical studies; we first show the values of the regularization coefficients  $a_{jk}$  by varying masks; then we show the input gradient values that are regularized by  $a_{jk}$  (9) after the MSDA training; finally, we show that the best choice of the mask design can be varying by the target task settings. In addition, we propose two generalized versions of Mixup and CutMix, called HMix and GMix, that empirically show the intermediate property of Mixup and CutMix.

**Introduction to HMix and GMix.** Recall that the regularization coefficients  $a_{jk}$  is determined by M (Equation (9)). For example, by choosing  $M = \lambda \vec{1}$  (*i.e.*, Mixup),  $a_{jk}$  is always  $(1 - \lambda)^2$ . On the other hand, the result slightly changes for CutMix:  $a_{jk}$  depends on how j and k are close. Informally, due to dependency between  $M_j$  and  $M_k$  (as M's component is always 0 in the cropped box regions and 1 in others), close j and k give large  $a_{jk}$ , but distant j and k give small  $a_{jk}$ .  $a_{jk}$  is calculated as

$$a_{jk} = \frac{\max(\min(h(j_1) - l(k_1), h(k_1) - l(j_1)), 0) \max(\min(h(j_2) - l(k_2), h(k_2) - l(j_2)), 0)}{(n - [\sqrt{1 - \lambda}n])^2}$$
(10)

where  $j = (j_1, j_2), k = (k_1, k_2), h(t) = \min(t, n - [\sqrt{1 - \lambda}n]), l(t) = \max(t - [\sqrt{1 - \lambda}n], 0).$ 

We visualize  $a_{jk}$  of different MSDA methods in Figure 4. We compare Mixup, CutMix, Stochastic Mixup & CutMix. We also propose two generalized MSDA methods, named **HMix** and **GMix**. Before comparing the methods, we first formally define Stochastic Mixup & CutMix, HMix and GMix. These methods can be formed as (1) where the definition of M is varying by the methods.

Stochastic Mixup & CutMix is a practical variant of MSDA by considering Mixup and CutMix at the same time. By a simple alternation of two augmentations, the state-of-the-art performances on large-scale datasets are shown [63, 68]. Stochastic Mixup & CutMix is the same as Equation (1) by setting  $M(\lambda) = (1 - \lambda)\vec{1}$  with probability q and  $M(\lambda) = M^{\text{cutmix}}(\lambda)$  with probability 1 - q. We choose q = 0.5 as [63, 68]. In our loss function perspective, the regularizing coefficient terms (*i.e.*,  $\mathcal{R}_2, \mathcal{R}_3$ ) become the average of Mixup and CutMix's regularization coefficient. Namely, let  $a_{ij}^{\text{mixup}} = (1 - \lambda)^2$  be a regularization coefficient of Mixup and  $a_{ij}^{\text{cutmix}}$  be a regularization coefficient of CutMix (10), then the regularization coefficients of Stochastic Mixup & CutMix is  $qa_{ij}^{\text{cutmix}} + (1 - q)a_{ij}^{\text{mixup}}$ .

Here, we additionally propose two MSDA variants, HMix and GMix, that leverage the advantages of Mixup and CutMix, resulting in showing the intermediate property between Mixup and CutMix.

*Hybrid version of Mixup and CutMix (HMix)* combines Mixup and CutMix by shrinking the CutMix cropped box region and linearly interpolating two images in the areas out of the box as Mixup. The shrinking ratio of the cropped box region is determined by the ratio r. HMix can be written as (1) by setting M by (1) randomly cropped box region with side length  $\sqrt{1 - \lambda}\sqrt{rN}$  where N is the side length of the original image, and make M's component in the box region as 0 (2) in the areas other than the box, we set  $M_i$  as  $\frac{\lambda}{1-(1-\lambda)r}$ . We can easily check that  $\mathbb{E}[M] = \lambda \vec{1}$ . As  $r \to 0$ , this



Figure 4: Visualization of regularization coefficients for different MSDA methods.  $a_{ij}$  values of Mixup, CutMix, Stochastic Mixup & CutMix (the alternation of Mixup and CutMix), HMix, GMix (described in Section 4 and Appendix H) are shown. Each (x, y) value is computed by  $\mathbb{E}_i a_{i,i+(x,y)}$  where *i* is a pixel vector.

method goes to Mixup, and as  $r \to 1$ , this method goes to CutMix. Note that the ratio r can be a random variable, such as  $Beta(\gamma, \gamma)$ . In this case, if we set  $\gamma \to 0$ , as  $Beta(\gamma, \gamma)$  goes to Bernoulli distribution, this is equivalent to Stochastic Mixup & CutMix.

We propose *Gaussian Mixup (GMix)* to relax the CutMix box condition to a continuous version as the rectangle cropping of CutMix causes implausible augmented data, *e.g.*, the boundary between two mixed samples. Therefore, we combine two ideas of Mixup and CutMix. Firstly, we select a point *p* from the given input. Then, we make  $M_i$  as the related function with  $||i - p||^2$ . Specifically, we use the Gaussian function for making M: (1) randomly select a point *p* in image (2) in the areas other than the box, we set  $M_i$  as  $1 - \exp\left(-\frac{||i-p||^2\pi}{2(1-\lambda)N^2}\right)$ . The proposed GMix has the following  $a_{ij}$ 

$$a_{ij} = \frac{1}{N^2} \sum_{p \in \text{pixel}} \exp\left(\frac{-\pi}{2(1-\lambda)N^2} \left(-\|i-p\|^2 - \|j-p\|^2\right)\right)$$
  
=  $\int_{\mathbb{R}^2} \exp\left(\frac{-\pi}{2(1-\lambda)N^2} \left(-\|i-p\|^2 - \|j-p\|^2\right)\right) dx$   
=  $(1-\lambda) \exp\left(\frac{-\pi}{(1-\lambda)N^2} \left\|\frac{i-j}{2}\right\|^2\right).$  (11)

As seen in Equation (11),  $a_{ij}$  smoothly goes down when the pixel distance becomes larger.

Figure 3 shows the examples generated Mixup, CutMix, HMix, and GMix. The proposed methods (HMix and GMix) generate images in a hybrid form with the properties of both Mixup and CutMix.

**Comparison in terms of regularization coefficients**  $a_{jk}$ . We illustrate the regularization coefficients  $a_{jk}$  of the different MSDA methods in Figure 4. In particular, we fix the mask parameter  $\lambda$  to 0.5 and the input resolution to 64 × 64. Figure 4 shows the difference between the MSDA methods in terms of how they regularize the input gradients and input Hessians: Mixup has equal weights to every gradient component or Hessian component, while CutMix gives high regularization in close coordinate gradient products or Hessian. We also observe that the hybrid methods (*e.g.*, Stochastic Mixup & CutMix, HMix, and GMix) show the intermediate coefficient values of Mixup and CutMix.

**Comparison in terms of the regularized input gradients after MSDA training.** Equation (8) shows that the regularization term  $a_{ij}$  directly affects to the pixel gradients  $|\partial_i f_\theta(x_k) \partial_j f_\theta(x_k)|$  in our approximated loss function. The purpose of Figure 5 is to show how the pixel gradients are actually regularized after training. We investigate the amount of the regularized input gradients by  $|\partial_v f_\theta(x) \partial_{v+p} f_\theta(x)|$  with respect to the pixel distance vector p for trained models by different MSDA methods. Here, if our approximated loss function actually behaves as a regularization, then we can expect that the pixel gradients  $|\partial_v f_\theta(x) \partial_{v+p} f_\theta(x)|$  is small when  $a_{ij}$  is large for the given p.

We first define the partial gradient product as follows:

$$\operatorname{PartialGradProd}(x, p) = \max_{v} \left| \partial_{v} f_{\theta}(x) \partial_{v+p} f_{\theta}(x) \right|$$
(12)

Now, we visualize the pixel-wise maximum values of PartialGradProd(x, p) in Figure 5. We train different models  $f_{\theta}$  on resized ImageNet (64 x 64) and measure the values on the validation dataset. The x-axis and y-axis of Figure 5 denote the pixel distance p along each x and y axis, and the scale



Figure 5: **Regularized input gradients by MSDA.** The normalized pixel-wise partial gradient norm product comparison among the models trained with vanilla setting (a), Mixup (b) and CutMix (c). We plotted (12), and x and y axis denote the pixel distance p along each axis.

Table 1: Different tasks need different MSDA strategies. Validation accuracies of Mixup and CutMix trained networks on two different scenarios on ImageNet-100. Each scenario assumes different pixel importances.

	Mixup	CutMix	$\Delta$ (CutMix - Mixup)
Scenario 1: Large crop	58.3	<b>64.4</b>	+6.1
Scenario 2: Small crop	<b>67.7</b>	67.0	-0.7

of the colorbar denotes the value of the maximum partial gradient product. In the figure, we can observe that CutMix reasonably regularizes effectively in the input gradients products when a pixel distance is small; these results aligned with our previous interpretation, CutMix behaves a pixel-level regularizer where it gives stronger regularization (larger  $a_{ij}$ ) to the closer pixels. Note that we are not discussing the relationship between regularizing effects and accuracy but discussing regularizing coefficients and the optimized function's pixel gradients.

Understanding application cases when a specific MSDA design choice works better than others. From our theoretical results and empirical studies, we have shown that the design choice of MSDAs (i.e., M) leads to different regularization effects by regularization coefficient  $a_{jk}$ . Furthermore, as we have shown in Theorem 2, there always exists a mask that can form any desired  $a_{jk}$ . We hypothesize that for the given dataset, if a short distance relation is relatively more important than longer distance relations, then CutMix will be better than Mixup. On the contrary, in the opposite case, if a short distance relation is relatively less important, then Mixup will be better than CutMix.

Here, we study different task scenarios when different  $a_{jk}$ s are required by controlling the pixel-level importance of ImageNet-100 [60] training images. In particular, we design two different scenarios where each of them needs different regularization strategies due to the different pixel-level importance of each task. The results are shown in Table 1. We also report the performance of our proposed methods in both scenarios 1 and 2 in Appendix G.

Scenario 1: Smaller objects by large crop size. We randomly crop a large region (80% to 100%) of an image and resize to  $64 \times 64$  to train a model. As the objects in the image become small, a close-distance relationship might be more important than a large-distance relationship. Here, we expect CutMix performs better than Mixup as shown in Table 1.

Scenario 2: Larger objects by small crop size We randomly crop a small region (25% to 40%) of an image and resize to  $64 \times 64$  to train a model. Contrary to Scenario 1, the objects in the image would become large in the cropping region and the large-distance relationship might be important, therefore, we expect that Mixup performs better than CutMix. This hypothesis is aligned to Table 1.

# 5 Comparison of Different MSDA Design Choices: An Empirical Validation

In this section, we compare various MSDA methods on two popular large-scale image classification benchmarks: CIFAR-100 [41] and ImageNet-1K [17]. We will confirm that our proposed design choices, HMix and GMix, are not only theoretically interpolating Mixup and CutMix in the toy

	_				
Augmentation Method	RN56	WRN28-2	PreActRN18	PreActRN34	PreActRN50
Vanilla (no MDSA)	73.23	73.50	76.73	77.68	79.07
Mixup	73.12	74.05	77.21	79.02	79.34
CutMix	74.83	74.79	78.66	80.05	81.23
PuzzleMix	-	76.51	79.38	80.89	82.46
Stochastic Mixup & CutMix	74.88	75.49	79.25	81.05	81.21
HMix (ours)	74.99	75.68	<b>79.25</b>	<b>81.07</b>	81.38 81.45
Givinx (ours)	15.15	/0.15	/9.1/	00.52	01.45

Table 2: CIFAR-100 classification. Comparison of various MSDA methods on various network architectures. Note that PuzzleMix needs additional computations (twice than others) for computing the input saliency.

settings, but also taking benefits of each method by showing great performances in real-world applications. The implementation details and the hyper-parameter study can be found in Appendix F.

Results on CIFAR-100 classification. We evaluate our method (HMix and GMix) against baseline MSDA methods including Mixup [74], CutMix [70], Stochastic Mixup & CutMix, [63] and PuzzleMix [39] on CIFAR-100 dataset. Here, we include PuzzleMix, to see the effectiveness of our data-agnostic method against the data-aware mask strategy. Note that although our theoretical results (Section 3) are based on the data-agnostic mask selection methods, our theoretical results can be easily extended to the data-dependent mask selection methods. We leave the extension as a future research direction.

To see the generalizability of our methods, we train various network architectures including ResNet-56 (RN56) [26], WideResNet28-2 (WRN28-2) [73], PreActResNet-18 (PreActRN18) [27], PreActResNet-34 (PreActRN34) [27] and PreActResNet-50 (PreActRN50) [27] with various MSDA methods. We train networks for 300 epochs using SGD optimizer with a learning rate 0.2. Table 2 shows the summarized results. We set the hyper-parameter  $\alpha$  for Mixup, CutMix, and Stochastic Mixup & CutMix to 1.  $\alpha$  for HMix and GMix were set to 1 and 0.5, respectively. We use r = 0.5for HMix. In the table, HMix and GMix outperform Mixup only and CutMix only counterparts and Stochastic Mixup & CutMix often show comparable performances to HMix and GMix. Our methods show comparable performance with the state-of-the-art data-dependent strategy PuzzleMix.

Results on ImageNet-1K classification. Ta- Table 3: ImageNet-1K classification. Comparison of ble 3 shows the comparison of various MSDA various MSDA methods on ResNet-50 architecture. methods on ImageNet-1K. We train ResNet-50 [26] with various MSDA methods for 300 epochs using SGD optimizer with a learning rate 0.1. We set the hyper-parameter  $\alpha$  for all methods except Mixup to 1, while Mixup has  $\alpha = 0.8$ . We use r = 0.75 for HMix. Here, we do not include PuzzleMix because it needs heavy additional computations to compute the input saliencies. In the table, HMix shows the

Augmentation Method	Top-1 accuracy
Vanilla (no MDSA)	75.68 (+0.00)
Mixup	77.78 (+2.10)
CutMix	78.04 (+2.36)
Stochastic Mixup & CutMix	78.13 (+2.45)
HMix (ours)	<b>78.38</b> (+2.70)
GMix (ours)	78.13 (+2.45)

best performance, while GMix and Stochastic Mixup & CutMix show the second-best performances. Evaluations on various robustness benchmarks are in Appendix G.

#### 6 Conclusion

We analyze MSDA by a unified theoretical framework. Our unified theoretical results show that any MSDA method behaves as a regularization on the input gradients and Hessians, where the degree of the regularization is controlled by the design choice of MSDA. We compare various MSDA methods in (1) regularization coefficient (2) regularized gradients (3) model performances in various scenarios with different pixel-level importance. We propose two simple MSDA methods, HMix and GMix, which leverage the benefits of Mixup and CutMix by their design. Our experimental results show that HMix and GMix outperform popular MSDA methods Mixup and CutMix. Furthermore, our methods

show comparable or outperformed performances than the state-of-the-art MSDA method, Stochastic Mixup & CutMix, in CIFAR-100 and ImageNet classification tasks.

# References

- R. Arora, P. Bartlett, P. Mianjy, and N. Srebro. Dropout: Explicit forms and capacity control. *ICML*, 2021. 22
- [2] P. L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002. 23, 24
- [3] C. Beckham, S. Honari, V. Verma, A. M. Lamb, F. Ghadiri, R. D. Hjelm, Y. Bengio, and C. Pal. On adversarial mixup resynthesis. *NeurIPS*, 2019. 1
- [4] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel. Mixmatch: A holistic approach to semi-supervised learning. *NeurIPS*, 2019. 1
- [5] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao. Yolov4: Optimal speed and accuracy of object detection. arXiv preprint arXiv:2004.10934, 2020. 2
- [6] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020. 1
- [7] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, et al. Api design for machine learning software: experiences from the scikit-learn project. arXiv preprint arXiv:1309.0238, 2013. 4
- [8] L. Carratino, M. Cissé, R. Jenatton, and J.-P. Vert. On mixup regularization. *arXiv preprint arXiv:2006.06049*, 2020. 2, 4
- [9] J. Cha, S. Chun, K. Lee, H.-C. Cho, S. Park, Y. Lee, and S. Park. Swad: Domain generalization by seeking flat minima. In *Neural Information Processing Systems (NeurIPS)*, 2021. 5
- [10] Y.-T. Chang, Q. Wang, W.-C. Hung, R. Piramuthu, Y.-H. Tsai, and M.-H. Yang. Mixupcam: Weakly-supervised semantic segmentation via uncertainty regularization. arXiv preprint arXiv:2008.01201, 2020. 1
- [11] P. Chen, S. Liu, H. Zhao, and J. Jia. Gridmask data augmentation. *arXiv preprint* arXiv:2001.04086, 2020. 1
- [12] M. Chidambaram, X. Wang, Y. Hu, C. Wu, and R. Ge. Towards understanding the data dependency of mixup-style training. *ICLR*, 2022. 2, 6, 25
- [13] S. Chun, S. J. Oh, S. Yun, D. Han, J. Choe, and Y. Yoo. An empirical evaluation on robustness and uncertainty of regularization methods. *ICML Workshop on Uncertainty and Robustness in Deep Learning*, 2019. 27
- [14] S. Chun and S. Park. Styleaugment: Learning texture de-biased representations by style augmentation without pre-defined textures. arXiv preprint arXiv:2108.10549, 2021. 1
- [15] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. Randaugment: Practical automated data augmentation with a reduced search space. CVPR, 2020. 1
- [16] A. Dabouei, S. Soleymani, F. Taherkhani, and N. M. Nasrabadi. Supervising the mixing data augmentation. CVPR, 2021. 1
- [17] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. CVPR, 2009. 9
- [18] M. Faramarzi, M. Amini, A. Badrinaaraayanan, V. Verma, and S. Chandar. Patchup: A regularization technique for convolutional neural networks. *arXiv preprint arXiv:2006.07794*, 2020. 1
- [19] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021. 5

- [20] G. French, S. Laine, T. Aila, M. Mackiewicz, and G. Finlayson. Semi-supervised semantic segmentation needs strong, varied perturbations. *BMVC*, 2020. 1
- [21] T. Garipov, P. Izmailov, D. Podoprikhin, D. P. Vetrov, and A. G. Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. In *Neural Information Processing Systems*, 2018. 5
- [22] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *ICLR*, 2015. 27
- [23] K. Greenewald, A. Gu, M. Yurochkin, J. Solomon, and E. Chien. k-mixup regularization for deep learning via optimal transport. arXiv preprint arXiv:2106.02933, 2021. 1, 3
- [24] H. Guo, Y. Mao, and R. Zhang. Augmenting data with mixup for sentence classification: An empirical study. arXiv preprint arXiv:1905.08941, 2019.
- [25] E. Harris, A. Marcu, M. Painter, M. Niranjan, A. Prügel-Bennett, and J. Hare. Fmix: Enhancing mixed sample data augmentation. arXiv preprint arXiv:2002.12047, 2020. 1, 28
- [26] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. CVPR, 2016. 10, 26
- [27] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. ECCV, 2016. 10
- [28] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *ICLR*, 2018. 27
- [29] M. Hong, J. Choi, and G. Kim. Stylemix: Separating content and style for enhanced data augmentation. CVPR, 2021. 28
- [30] H. Inoue. Data augmentation by pairing samples for images classification. *arXiv preprint arXiv:1801.02929*, 2018. 1
- [31] P. Izmailov, D. Podoprikhin, T. Garipov, D. Vetrov, and A. G. Wilson. Averaging weights leads to wider optima and better generalization. *Conference on Uncertainty in Artificial Intelligence*, 2018. 5
- [32] J. Jeong, S. Cha, Y. Yoo, S. Yun, T. Moon, and J. Choi. Observations on k-image expansion of image-mixing augmentation for classification. arXiv preprint arXiv:2110.04248, 2021. 1, 3
- [33] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *ICML*, 2021. 1
- [34] A. Jindal, D. Gnaneshwar, R. Sawhney, and R. R. Shah. Leveraging bert with mixup for sentence classification (student abstract). AAAI, 2020. 1
- [35] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus. Hard negative mixing for contrastive learning. *NeurIPS*, 2020. 1
- [36] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference* on Learning Representations, 2017. 5
- [37] G. Kim, D. K. Han, and H. Ko. Specmix: A mixed sample data augmentation method for training withtime-frequency domain features. *INTERSPEECH*, 2021. 1
- [38] J. H. Kim, W. Choo, H. Jeong, and H. O. Song. Co-mixup: Saliency guided joint mixup with supermodular diversity. *ICLR*, 2021. 1, 3, 5, 28
- [39] J. H. Kim, W. Choo, and H. O. Song. Puzzle mix: Exploiting saliency and local statistics for optimal mixup. *ICML*, 2020. 1, 3, 5, 10, 26, 28
- [40] S. Kim, G. Lee, S. Bae, and S.-Y. Yun. Mixco: Mix-up contrastive learning for visual representation. arXiv preprint arXiv:2010.06300, 2020. 1

- [41] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images, 2009. 9
- [42] K. Lee, Y. Zhu, K. Sohn, C.-L. Li, J. Shin, and H. Lee. \$i\$-mix: A domain-agnostic strategy for contrastive representation learning. *ICLR*, 2021. 1, 2
- [43] J. Li, R. Socher, and S. C. Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. *ICLR*, 2020. 1
- [44] Z. Liu, S. Li, G. Wang, C. Tan, L. Wu, and S. Z. Li. Decoupled mixup for data-efficient learning. arXiv preprint arXiv:2203.10761, 2022. 1
- [45] Z. Liu, S. Li, D. Wu, Z. Chen, L. Wu, J. Guo, and S. Z. Li. Automix: Unveiling the power of mixup. ECCV, 2022. 1, 28
- [46] C. Ma and L. Ying. On linear stability of sgd and input-smoothness of neural networks. *NeurIPS*, 2021. 5, 6
- [47] D. Mahajan, R. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, and L. Van Der Maaten. Exploring the limits of weakly supervised pretraining. *ECCV*, 2018. 1
- [48] I. Medennikov, Y. Y. Khokhlov, A. Romanenko, D. Popov, N. A. Tomashenko, I. Sorokin, and A. Zatvornitskiy. An investigation of mixup training strategies for acoustic models in asr. *Interspeech*, 2018. 1
- [49] L. Meng, J. Xu, X. Tan, J. Wang, T. Qin, and B. Xu. Mixspeech: Data augmentation for low-resource automatic speech recognition. *ICASSP*, 2021. 1
- [50] D. Misra. Mish: A self regularized non-monotonic activation function. BMVC, 2020. 2
- [51] S.-M. Moosavi-Dezfooli, A. Fawzi, J. Uesato, and P. Frossard. Robustness via curvature regularization, and vice versa. *CVPR*, 2019. 6
- [52] W. Mustafa, R. A. Vandermeulen, and M. Kloft. Input hessian regularization of neural networks. arXiv preprint arXiv:2009.06571, 2020. 6
- [53] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le. Specaugment: A simple data augmentation method for automatic speech recognition. In *Interspeech*, 2019. 1
- [54] J. Qin, J. Fang, Q. Zhang, W. Liu, X. Wang, and X. Wang. Resizemix: Mixing data with preserved object information and true labels. arXiv preprint arXiv:2012.11101, 2020. 1, 28
- [55] S.-A. Rebuffi, S. Gowal, D. A. Calian, F. Stimberg, O. Wiles, and T. A. Mann. Data augmentation can improve robustness. *NeurIPS*, 2021. 21
- [56] S. Ren, H. Wang, Z. Gao, S. He, A. Yuille, Y. Zhou, and C. Xie. A simple data mixing prior for improving self-supervised learning. *CVPR*, 2022. 2
- [57] J. Y. Sohn, L. Shang, H. Chen, J. Moon, D. Papailiopoulos, and K. Lee. Genlabel: Mixup relabeling using generative models. arXiv preprint arXiv:2201.02354, 2022. 1
- [58] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, and C.-L. Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *NeurIPS*, 2020. 1
- [59] R. Takahashi, T. Matsubara, and K. Uehara. Data augmentation using random image cropping and patching for deep cnns. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(9):2917–2931, 2019. 1, 3
- [60] Y. Tian, D. Krishnan, and P. Isola. Contrastive multiview coding. *arXiv preprint* arXiv:1906.05849, 2019. 9
- [61] Y. Tokozume, Y. Ushiku, and T. Harada. Between-class learning for image classification. *CVPR*, 2018. 1
- [62] Y. Tokozume, Y. Ushiku, and T. Harada. Learning from between-class examples for deep sound recognition. *ICLR*, 2018. 1

- [63] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient image transformers & distillation through attention. *ICML*, 2021. 1, 2, 7, 10
- [64] A. Uddin, M. Monira, W. Shin, T. Chung, S. H. Bae, et al. Saliencymix: A saliency guided data augmentation strategy for better regularization. *ICLR*, 2021. 1, 5, 28
- [65] V. Verma, A. Lamb, C. Beckham, A. Najafi, I. Mitliagkas, D. Lopez-Paz, and Y. Bengio. Manifold mixup: Better representations by interpolating hidden states. *ICML*, 2019. 1, 28
- [66] S. Wager, S. Wang, and P. S. Liang. Dropout training as adaptive regularization. *NeurIPS*, 2013.
   4
- [67] D. Walawalkar, Z. Shen, Z. Liu, and M. Savvides. Attentive cutmix: An enhanced data augmentation approach for deep learning based image classification. arXiv preprint arXiv:2003.13048, 2020. 1, 3
- [68] R. Wightman, H. Touvron, and H. Jégou. Resnet strikes back: An improved training procedure in timm. *NeurIPS*, 2021. 7
- [69] H. Yao, L.-K. Huang, L. Zhang, Y. Wei, L. Tian, J. Zou, J. Huang, et al. Improving generalization in meta-learning via task augmentation. *ICML*, 2021. 1, 2
- [70] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. *ICCV*, 2019. 1, 2, 3, 6, 10, 21, 27
- [71] S. Yun, S. J. Oh, B. Heo, D. Han, J. Choe, and S. Chun. Re-labeling imagenet: from single to multi-labels, from global to localized labels. *CVPR*, 2021. 2
- [72] S. Yun, S. J. Oh, B. Heo, D. Han, and J. Kim. Videomix: Rethinking data augmentation for video classification. arXiv preprint arXiv:2012.03457, 2020. 1
- [73] S. Zagoruyko and N. Komodakis. Wide residual networks. *BMVC*, 2016. 10
- [74] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. *ICLR*, 2018. 1, 2, 3, 6, 10, 21
- [75] L. Zhang, Z. Deng, K. Kawaguchi, A. Ghorbani, and J. Zou. How does mixup help with robustness and generalization? *ICLR*, 2021. 2, 3, 4, 5, 18, 20, 21, 22
- [76] L. Zhang, Z. Deng, K. Kawaguchi, and J. Zou. When and how mixup improves calibration. arXiv preprint arXiv:2102.06289, 2021. 2

# Appendix

We include additional materials in this document, including additional theoretical results (Appendix A, B, C, D, E), experimental details (Appendix F), additional experiments (Appendix G), and additional explanation for various MSDAs based on our analysis (Appendix H, I).

# **Negative Societal Impacts**

Our work may help the situation that data is insufficient, as other data augmentation papers do. We mainly proposed the theoretical aspect of MSDA, so there are no negative societal impacts.

# A Proof of Theorem 1

For convenience, we restate Theorem 1.

**Theorem 1.** Consider a loss function  $l \in \mathcal{L}$ . We define  $\tilde{D}_{\lambda}$  as  $\frac{\alpha}{\alpha+\beta}$ Beta $(\alpha+1,\beta)+\frac{\beta}{\alpha+\beta}$ Beta $(\beta+1,\alpha)$ . Assume that  $\mathbb{E}_{r_x \sim \mathcal{D}_X}[r_x] = 0$ . Then, we can re-write the general MSDA loss (4) as

$$L_m^{\text{MSDA}}(\theta) = L_m(\theta) + \sum_{i=1}^3 \mathcal{R}_i^{(\text{MSDA})}(\theta) + \mathbb{E}_{\lambda \sim \tilde{\mathcal{D}}(\lambda)} \mathbb{E}_M[(1-M)^{\mathsf{T}} \varphi(1-M)(1-M)],$$

where  $\lim_{a\to 0} \varphi(a) = 0$ ,

$$\mathcal{R}_{1}^{(\text{MSDA})}(\theta) = \frac{1}{m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}} (\nabla f_{\theta}(x_{i}) \odot (r_{x} - x_{i}))^{\mathsf{T}} \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbb{E}_{M} (1 - M)$$
$$= \frac{1}{m} \sum_{i=1}^{m} (y_{i} - h'(f_{\theta}(x_{i}))) (\nabla f_{\theta}(x_{i})^{\mathsf{T}} x_{i}) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} (1 - \lambda),$$
$$\mathcal{R}_{2}^{(\text{MSDA})}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbf{G}(\mathcal{D}_{X}, x_{i}, f, M),$$
$$\mathcal{R}_{3}^{(\text{MSDA})}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbf{H}(\mathcal{D}_{X}, x_{i}, f, M),$$

and

$$\mathbf{G}(\mathcal{D}_X, x_i, f, M) = \mathbb{E}_M (1 - M)^{\mathsf{T}} \mathbb{E}_{r_x \sim \mathcal{D}_X} \left( \nabla f(x_i) \odot (r_x - x_i) \left( \nabla f(x_i) \odot (r_x - x_i) \right)^{\mathsf{T}} \right) (1 - M)$$

$$= \sum_{j,k \in \text{coord}} a_{jk} \partial_j f_\theta(x_i) \partial_k f_\theta(x_i) \left( \mathbb{E}_{r_x \sim \mathcal{D}_X} [r_{xj} r_{xk}] + x_{ij} x_{ik} \right),$$

$$\mathbf{H}(\mathcal{D}_X, x_i, f, M) = \mathbb{E}_{r_x \sim \mathcal{D}_X} \mathbb{E}_M (1 - M)^{\mathsf{T}} \left( \nabla^2 f_\theta(x_i) \odot \left( (r_x - x_i) (r_x - x_i)^{\mathsf{T}} \right) \right) (1 - M)$$

$$= \sum_{j,k \in \text{coord}} a_{jk} \left( \mathbb{E}_{r_x \sim \mathcal{D}_X} [r_{xj} r_{xk} \partial_{jk}^2 f_\theta(x_i)] + x_{ij} x_{ik} \partial_{jk}^2 f_\theta(x_i) \right),$$

where

$$a_{jk} := \mathbb{E}_M[(1 - M_j)(1 - M_k)].$$

*Proof of Theorem 1.* Due to the assumption of the theorem, we can rewrite the empirical loss for the non-augmented population as

$$L_m(\theta) = \frac{1}{m} \sum_{i=1}^m l(\theta, z_i) = \frac{1}{m} \sum_{i=1}^n [h(f_\theta(x_i)) - y_i f_\theta(x_i)].$$

Similarly, we can rewrite the MSDA loss as

$$\begin{split} L_m^{\text{MSDA}}(\theta) &= \frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E}_{\lambda \sim \mathcal{D}_\lambda} \mathbb{E}_M l(\theta, \tilde{z}_{i,j}^{(\text{MSDA})}(\lambda, 1 - \lambda)) \\ &= \frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E}_{\lambda \sim \text{Beta}(\alpha,\beta)} \mathbb{E}_M [h(f_\theta(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) - \tilde{y}_{i,j}^{(\text{MSDA})}(\lambda, 1 - \lambda) f_\theta(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))]. \end{split}$$

Putting the definition of  $\tilde{z}_{i,j}^{(\mathrm{MSDA})}(\lambda,1-\lambda)$  to the equation above, we have

$$\begin{split} L_{m}^{\text{MSDA}}(\theta) &= \frac{1}{m^{2}} \sum_{i,j=1}^{m} \Biggl( \mathbb{E}_{\lambda \sim \text{Beta}(\alpha,\beta)} \mathbb{E}_{M} \Biggl( \lambda(h(f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) - y_{i}f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(\lambda, 1 - \lambda))) \\ &+ (1 - \lambda)(h(f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) - y_{j}f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) \Biggr) \Biggr) \Biggr) \\ &= \frac{1}{m^{2}} \sum_{i,j=1}^{m} \Biggl( \mathbb{E}_{\lambda \sim \text{Beta}(\alpha,\beta)} \mathbb{E}_{B \sim \text{Bin}(\lambda)} \mathbb{E}_{M} \Biggl( B(h(f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) - y_{i}f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) - y_{i}f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) \Biggr) \\ &+ (1 - B)(h(f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) - y_{j}f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(M, 1 - M))) \Biggr) \Biggr) \Biggr). \end{split}$$

Note that  $\lambda \sim \text{Beta}(\alpha, \beta)$  and  $B|\lambda \sim \text{Bin}(\lambda)$ . By conjugacy, we can write the joint distribution of  $(\lambda, B)$  as

$$B \sim \operatorname{Bin}\left(\frac{\alpha}{\alpha+\beta}\right), \qquad \lambda|B \sim \operatorname{Beta}(\alpha+B,\beta+1-B).$$

Therefore, we have

$$L_{m}^{\text{MSDA}}(\theta) = \frac{1}{m^{2}} \sum_{i,j=1}^{m} \left( \mathbb{E}_{\lambda \sim \text{Beta}(\alpha+1,\beta)} \mathbb{E}_{M} \frac{\alpha}{\alpha+\beta} \left( h(f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda))) - y_{i}f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda)) \right) \right. \\ \left. + \mathbb{E}_{\lambda \sim \text{Beta}(\alpha,\beta+1)} \mathbb{E}_{M} \frac{\beta}{\alpha+\beta} \left( h(f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda))) - y_{j}f_{\theta}(\tilde{x}_{i,j}^{(\text{MSDA})}(\lambda, 1-\lambda)) \right) \right) \\ \left. = \frac{1}{m} \sum_{i=1}^{n} \mathbb{E}_{\lambda \sim \tilde{\mathcal{D}}(\lambda)} \mathbb{E}_{r_{x} \sim \mathcal{D}_{x}} \mathbb{E}_{M} \left[ h\left( f_{\theta}(M \odot x_{i} + (1-M) \odot r_{x}) \right) - y_{i}f_{\theta}(M \odot x_{i} + (1-M) \odot r_{x}) \right) \right] \right]$$
(13)

$$= \frac{1}{m} \sum_{i=1}^{n} \mathbb{E}_{\lambda \sim \tilde{\mathcal{D}}(\lambda)} \mathbb{E}_{r_x \sim \mathcal{D}_x} \mathbb{E}_M l(\theta, \hat{z}_i),$$
(14)

where  $\hat{z}_i = (M \odot x_i + (1 - M) \odot r_x, y_i).$ 

Let N = 1 - M. By defining  $\phi_i(N) = h \left( f_\theta(x_i + N \odot (r_x - x_i)) \right) - y_i f_\theta(x_i + N \odot (r_x - x_i))$ and applying Taylor expansion, we have

$$\phi_i(N) = \phi_i(0) + \nabla_N \phi_i(0)^{\mathsf{T}} N + \frac{1}{2} N^{\mathsf{T}} \nabla_N^2 \phi_i(0) N + N^{\mathsf{T}} \varphi(N) N,$$
(15)

where  $\lim_{N\to 0}\varphi(N)=0.$  Firstly, we calculate  $\phi_i(0)$  by

$$\phi_i(0) = h(f_\theta(x_i)) - y_i f_\theta(x_i). \tag{16}$$

Second, we calculate  $\nabla_N \phi_i(0)$  by

$$\frac{\partial \phi_i(N)}{\partial N_k} = \left(h'\left(f_\theta\left(x_i + N \odot (r_x - x_i)\right)\right) - y_i\right) \frac{\partial f_\theta}{\partial x_{ik}} \left(x_i + N \odot (r_x - x_i)\right) (r_{xk} - x_{ik}),$$

where we denote  $N_k$  as the kth element of N,  $x_{ik}$  as the kth element of  $x_i$ , and  $r_{xk}$  as the kth element of  $r_x$ . Therefore, we have

$$\nabla_N \phi_i(0)^{\mathsf{T}} N = (h'(f_\theta(x_i)) - y_i) \sum_k \left( \frac{\partial f_\theta}{\partial x_{ik}}(x_i)(r_{xk} - x_{ik}) \right) N_k$$
$$= (h'(f_\theta(x_i)) - y_i) \left( \nabla f \odot (r_x - x_i) \right) \cdot N.$$
(17)

Finally, we calculate  $\nabla^2_N \varphi_i(\vec{0})^T$  by

$$\begin{aligned} \frac{\partial^2 \phi_k(N)}{\partial N_k \partial N_j} &= \frac{\partial}{\partial N_j} \left( \left( h' \left( f_\theta \left( x_i + N \odot \left( r_x - x_i \right) \right) \right) - y_i \right) \frac{\partial f_\theta}{\partial x_{ik}} \left( x_i + N \odot \left( r_x - x_i \right) \right) \left( r_{xk} - x_{ik} \right) \right) \\ &= h'' \left( f_\theta \left( x_i + N \odot \left( r_x - x_i \right) \right) \right) \\ &\times \frac{\partial f_\theta}{\partial x_{ik}} \left( x_i + N \odot \left( r_x - x_i \right) \right) \left( r_{xk} - x_{ik} \right) \frac{\partial f_\theta}{\partial x_{ij}} \left( x_i + N \odot \left( r_x - x_i \right) \right) \left( r_{xj} - x_{ij} \right) \\ &+ \left( h' \left( f_\theta \left( x_i + N \odot \left( r_x - x_i \right) \right) \right) - y_i \right) \\ &\times \frac{\partial^2 f_\theta}{\partial x_{ik} \partial x_{ij}} \left( x_i + N \odot \left( r_x - x_i \right) \right) \left( r_{xk} - x_{ik} \right) (r_{xj} - x_{ij}). \end{aligned}$$

Therefore, we have

$$\frac{1}{2}N^{\mathsf{T}}\nabla_{N}^{2}\phi_{i}(0)N = \frac{1}{2}h^{\prime\prime}\left(f_{\theta}(x_{i})\right)\sum_{k,j}\left(\frac{\partial f_{\theta}}{\partial x_{ik}}(x_{i})(r_{xk}-x_{ik})\frac{\partial f_{\theta}}{\partial x_{ij}}(x_{i})(r_{xj}-x_{ij})N_{k}N_{j}\right) \\
+ \frac{1}{2}\left(h^{\prime}\left(f_{\theta}(x_{i})\right)-y_{i}\right)\sum_{k,j}\frac{\partial^{2}f_{\theta}}{\partial x_{ik}\partial x_{ij}}\left(x_{i}\right)\left(r_{xk}-x_{ik})(r_{xj}-x_{ij})N_{k}N_{j}\right) \\
= \frac{1}{2}h^{\prime\prime}\left(f_{\theta}(x_{i})\right)N^{\mathsf{T}}\left(\left(\nabla f\odot\left(r_{x}-x_{i}\right)\right)\left(\nabla f\odot\left(r_{x}-x_{i}\right)\right)^{\mathsf{T}}\right)N \\
+ \frac{1}{2}\left(h^{\prime}\left(f_{\theta}(x_{i})\right)-y_{i}\right)N^{\mathsf{T}}\left(\nabla^{2}f_{\theta}(x_{i})\odot\left((r_{x}-x_{i})(r_{x}-x_{i})^{\mathsf{T}}\right)\right)N. \tag{18}$$

Applying (16) - (18) to (15),

$$\phi_i(N) = (h(f_\theta(x_i)) - y_i f_\theta(x_i)) + (h'(f_\theta(x_i)) - y_i) (\nabla f \odot (r_x - x_i)) \cdot N + \frac{1}{2} h''(f_\theta(x_i)) N^{\mathsf{T}} ((\nabla f \odot (r_x - x_i)) (\nabla f \odot (r_x - x_i))^{\mathsf{T}}) N + \frac{1}{2} (h'(f_\theta(x_i)) - y_i) N^{\mathsf{T}} (\nabla^2 f_\theta(x_i) \odot ((r_x - x_i)(r_x - x_i)^{\mathsf{T}})) N + N^{\mathsf{T}} \varphi(N) N$$
(19)

Plugging (19) to (13), we conclude

$$L_m^{\text{MSDA}}(\theta) = \frac{1}{m} \sum_{i=1}^n \mathbb{E}_{\lambda \sim \tilde{\mathcal{D}}(\lambda)} \mathbb{E}_{r_x \sim \mathcal{D}_x} \mathbb{E}_M \phi(1-M)$$
  
=  $L_m(\theta) + \mathcal{R}_1(\theta) + \mathcal{R}_2(\theta) + \mathcal{R}_3(\theta) + \mathbb{E}_{\lambda \sim \tilde{\mathcal{D}}(\lambda)} \mathbb{E}_M[(1-M)^{\mathsf{T}} \varphi(1-M)M],$ 

where

$$\mathcal{R}_{1}(\theta) = \frac{1}{m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \left( \nabla f_{\theta}(x_{i}) \odot \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[r_{x} - x_{i}] \right) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbb{E}_{M} (1 - M),$$

$$\mathcal{R}_{2}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbb{E}_{M} (1 - M)^{\mathsf{T}} \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}} \left[ \nabla f(x_{i}) \odot (r_{x} - x_{i}) \left( \nabla f(x_{i}) \odot (r_{x} - x_{i}) \right)^{\mathsf{T}} \right] (1 - M),$$

$$\mathcal{R}_{3}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbb{E}_{M} (1 - M)^{\mathsf{T}} \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}} \left[ \nabla^{2} f_{\theta}(x_{i}) \odot ((r_{x} - x_{i})(r_{x} - x_{i})^{\mathsf{T}} \right] (1 - M).$$

# **B** Extension of Mixup: *n*-Mixup

In this section, due to notational complexity, we give the approximate loss function of the *n*-sample Mixup (*n*-Mixup). The same analysis can be applied to *n*-sample mixing strategy. We will define *n*-Mixup as followings. Mixup from the  $\mathbf{i} = (i_1, i_2, \ldots, i_n)$ th samples with  $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \ldots, \lambda_n)$  which is drawn from  $\mathcal{D}_{\Lambda}$  (mainly Dirichlet distribution), is defined as  $\tilde{z}_{\mathbf{i}} = \sum_{k=1}^{n} \lambda_k z_{i_k}$ . Similarly, we can define the *n*-Mixup loss as

$$L_m^{\text{n-mixup}}(\theta) = \mathbb{E}_{\mathbf{i} \sim \text{Unif}([m])} \mathbb{E}_{\boldsymbol{\lambda} \sim \mathcal{D}_{\Lambda}} l(\theta, \tilde{z}_{\mathbf{i}}(\boldsymbol{\lambda})) = \frac{1}{m^n} \sum_{\mathbf{i}} \mathbb{E}_{\boldsymbol{\lambda} \sim \mathcal{D}_{\Lambda}} l(\theta, \tilde{z}_{\mathbf{i}}(\boldsymbol{\lambda})).$$

Throughout this section, we consider  $\mathcal{D}_{\Lambda}$  as Dirichlet distribution (i.e.  $\mathcal{D}_{\Lambda} = \text{Dir}(\boldsymbol{\alpha}) = \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_n)$ ), which is the natural extension of Beta distribution.

**Theorem 5.** Consider the loss function in  $l \in \mathcal{L}$ . Then, we can rewrite the n-Mixup loss as

$$L_m^{n-mix}(\theta) = \frac{1}{m} \sum_{k=1}^n \mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}(\Lambda)} \mathbb{E}_{r_{x,2}, \cdots, r_{x,n} \sim \mathcal{D}_x} \varphi_k(\lambda_2, \cdots, \lambda_n)$$
$$= L_m(\theta) + \mathcal{R}_1(\theta) + \mathcal{R}_2(\theta) + \mathcal{R}_3(\theta) + \mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}(\Lambda)} [o(\|(\lambda_2, \cdots, \lambda_n)\|^2)],$$

where

$$\mathcal{R}_1(\theta) = \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[1-\lambda_1]}{m} \sum_{i=1}^m (h'(f_{\theta}(x_i)) - y_i) \nabla f_{\theta}(x_i)^T \mathbb{E}_{r_x \sim \mathcal{D}_X}[r_x - x_i],$$

$$\mathcal{R}_{2}(\theta) = \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[\sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \nabla f_{\theta}(x_{i})^{T} \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})(r_{x} - x_{i})^{T}] \nabla_{\theta} f(x_{i})$$
$$+ \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[(1 - \lambda_{1})^{2} - \sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \nabla f_{\theta}(x_{i})^{T} \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})] \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})^{T}] \nabla_{\theta} f(x_{i}),$$

$$\mathcal{R}_{3}(\theta) = \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[\sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})\nabla^{2}f_{\theta}(x_{i})(r_{x} - x_{i})^{T}] \\ + \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[(1 - \lambda_{1})^{2} - \sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})]\nabla^{2}f_{\theta}(x_{i}) \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})^{T}].$$

As Mixup's approximate loss function [75] or Theorem 1, *n*-Mixup also regularizes  $\nabla f$  and  $\nabla^2 f$ .

*Proof.* Due to the assumption of the theorem, we can rewrite the empirical loss for the non-augmented population as

$$L_m(\theta) = \frac{1}{m} \sum_{i=1}^m l(\theta, z_i) = \frac{1}{m} \sum_{i=1}^n [h(f_\theta(x_i)) - y_i f_\theta(x_i)].$$

Similarly, we can rewrite the n-Mixup loss as

$$L_m^{n-\min}(\theta) = \frac{1}{m^n} \sum_{\mathbf{i}} \mathbb{E}_{\boldsymbol{\lambda} \sim \mathcal{D}_{\Lambda}} l(\theta, \tilde{z}_{\mathbf{i}}(\boldsymbol{\lambda})) = \frac{1}{m^n} \mathbb{E}_{\boldsymbol{\lambda} \sim \text{Dir}(\boldsymbol{\alpha})} \sum_{\mathbf{i}} [h(f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda}))) - \tilde{y}_{\mathbf{i}}(\boldsymbol{\lambda}) f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda}))].$$

Putting the definition of  $\tilde{x}_i$  and  $\tilde{y}_i$  to the equation above, we have

$$L_m^{n-\min}(\theta) = \frac{1}{m^n} \mathbb{E}_{\boldsymbol{\lambda} \sim \text{Dir}(\boldsymbol{\alpha})} \sum_{\mathbf{i}} \sum_{k=1}^n \lambda_k (h(f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda}))) - y_k f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda})))$$
$$= \frac{1}{m^n} \sum_{\mathbf{i}} \mathbb{E}_{\boldsymbol{\lambda} \sim \text{Dir}(\boldsymbol{\alpha})} \mathbb{E}_{\boldsymbol{\beta} \sim \text{Mult}(\boldsymbol{\lambda})} \sum_{k=1}^n \beta_k (h(f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda}))) - y_k f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda}))),$$

where  $\boldsymbol{\beta} = (\beta_1, \beta_2, \cdots, \beta_n)$  and  $\text{Mult}(\boldsymbol{\lambda})$  is multinomial distribution. Note that  $\boldsymbol{\lambda} \sim \text{Dir}(\boldsymbol{\alpha})$  and  $\boldsymbol{\beta} | \boldsymbol{\lambda} \sim \text{Mult}(\boldsymbol{\lambda})$ . By conjugacy, we can write the joint distribution of  $(\boldsymbol{\alpha}, \boldsymbol{\beta})$  as

$$\boldsymbol{\beta} \sim \operatorname{Mult}\left(\frac{\boldsymbol{\alpha}}{\sum \alpha_i}\right), \qquad \boldsymbol{\lambda} | \boldsymbol{\beta} \sim \operatorname{Dir}(\boldsymbol{\alpha} + \boldsymbol{\beta}).$$

Therefore,

$$L_m^{n-\min}(\theta) = \frac{1}{m^n} \sum_{\mathbf{i}} \sum_{k=1}^n \frac{\alpha_k}{\sum \alpha_i} \mathbb{E}_{\boldsymbol{\lambda} \sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_k+1, \dots, \alpha_n)} (h(f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda}))) - y_k f_{\theta}(\tilde{x}_{\mathbf{i}}(\boldsymbol{\lambda}))) = \frac{1}{m} \sum_{k=1}^n \mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}(\Lambda)} \mathbb{E}_{r_{x,2}, \cdots, r_{x,n} \sim \mathcal{D}_x} \left[ h\left( f_{\theta}\left( \lambda_1 x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) \right) - y_k f_{\theta}\left( \lambda_1 x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) \right],$$

$$(20)$$

where  $\tilde{D}_{\Lambda} = \sum_{l=1}^{n} \frac{\alpha_{l}}{\sum \alpha_{i}} \operatorname{Dir}(\alpha_{l}+1, \alpha_{l+1}, \dots, \alpha_{l+n-1})$  and  $\mathcal{D}_{x}$  is the empirical distribution induced by training samples. We regard the index with mod n. Defining  $\varphi_{k}(\cdot)$  as

$$\varphi_k(\lambda_2, \cdots, \lambda_n) = h\left(f_\theta\left((1 - \sum_{j=2}^n \lambda_j)x_k + \sum_{j=2}^n \lambda_j r_{x,j}\right)\right) - y_k f_\theta\left((1 - \sum_{j=2}^n \lambda_j)x_k + \sum_{j=2}^n \lambda_j r_{x,j}\right),$$
  
we can use twice the differentiability of  $f(\cdot)$  and  $h(\cdot)$  so we have

we can use twice the differentiability of  $f(\cdot)$  and  $h(\cdot)$ , so we have

$$\varphi_k(\lambda_2,\cdots,\lambda_n) = \varphi_k(\vec{0}) + \nabla\varphi_k(\vec{0})^T(\lambda_2,\cdots,\lambda_n) + \frac{1}{2}(\lambda_2,\cdots,\lambda_n)^T \nabla^2 \varphi_k(\vec{0})^T(\lambda_2,\cdots,\lambda_n)^T + o(\|(\lambda_2,\cdots,\lambda_n)\|^2)$$
(21)

Firstly, we calculate  $\varphi_k(\vec{0})$  by

$$\varphi_k(\vec{0}) = h(f_\theta(x_k)) - y_k f_\theta(x_k).$$
(22)

Second, we calculate  $\nabla \varphi_k(\vec{0})$  by

$$\frac{\partial \varphi_k(\lambda_2, \cdots, \lambda_n)}{\partial \lambda_i} = h' \left( f_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) \right) f_\theta' \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) (r_{x,i} - x_k) - y_k f_\theta' \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) (r_{x,i} - x_k)$$

Therefore, we have

$$\frac{\partial \varphi_k(\lambda_2, \cdots, \lambda_n)}{\partial \lambda_i} \Big|_{(\lambda_2, \cdots, \lambda_n) = \vec{0}} = (h'(f_\theta(x_k)) - y_i) \nabla f_\theta(x_k)^T (r_{x,i} - x_k).$$
(23)

Finally, we calculate  $\nabla^2 \varphi_i(\vec{0})^T$  by

$$\begin{split} \frac{\partial^2 \varphi_k(\lambda_2, \cdots, \lambda_n)}{\partial \lambda_i \partial \lambda_s} &= \frac{\partial}{\partial \lambda_s} \left( h' \left( f_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) \right) f'_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) (r_{x,i} - x_k) \right) \\ &\quad - y_k f'_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) (r_{x,i} - x_k) \right) \\ &= h'' \left( f_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) \right) \left[ f'_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) (r_{x,i} - x_k) \right] \\ &\quad \left[ f'_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) (r_{x,s} - x_k) \right] \\ &\quad + h' \left( f_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) \right) (r_{x,i} - x_k)^T \nabla^2 f_\theta \left( (1 - \sum_{j=2}^n \lambda_j) x_k + \sum_{j=2}^n \lambda_j r_{x,j} \right) (r_{x,s} - x_k) . \end{split}$$

Therefore,

$$\frac{\partial^2 \varphi_k(\lambda_2, \cdots, \lambda_n)}{\partial \lambda_i \partial \lambda_s} \Big|_{(\lambda_2, \cdots, \lambda_n) = \vec{0}} = (h'' \left( f_\theta(x_k) - y_k \right)) \left( \nabla f_\theta \left( x_k \right)^T \left( r_{x,i} - x_k \right) \right) \left( \nabla f_\theta \left( x_k \right)^T \left( r_{x,s} - x_k \right) \right) - y_k (r_{x,i} - x_k)^T \nabla^2 f_\theta \left( x_k \right) (r_{x,s} - x_k).$$
(24)

Applying (22) - (24) to (21),

$$\varphi_{k}(\lambda_{2}, \cdots, \lambda_{n}) = (h(f_{\theta}(x_{k})) - y_{k}f_{\theta}(x_{k})) + \sum_{j=2}^{n} (h'(f_{\theta}(x_{k})) - y_{i})\nabla f_{\theta}(x_{k})^{T}(r_{x,i} - x_{k})\lambda_{j}$$
  
+  $\frac{1}{2}\sum_{i,s=2}^{n} \left( (h''(f_{\theta}(x_{k}) - y_{k})) \left( \nabla f_{\theta}(x_{k})^{T}(r_{x,i} - x_{k}) \right) \left( \nabla f_{\theta}(x_{k})^{T}(r_{x,s} - x_{k}) \right)$   
-  $y_{k}(r_{x,i} - x_{k})^{T} \nabla^{2} f_{\theta}(x_{k}) (r_{x,s} - x_{k}) \right) \lambda_{i}\lambda_{s} + o(\|(\lambda_{2}, \cdots, \lambda_{n})\|^{2}).$  (25)

Plugging (25) to (20),

$$L_m^{n-\min}(\theta) = \frac{1}{m} \sum_{k=1}^n \mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}(\Lambda)} \mathbb{E}_{r_{x,2}, \cdots, r_{x,n} \sim \mathcal{D}_x} \varphi_k(\lambda_2, \cdots, \lambda_n)$$
$$= L_m(\theta) + \mathcal{R}_1(\theta), + \mathcal{R}_2(\theta) + \mathcal{R}_3(\theta) + \mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}(\Lambda)}[o(\|(\lambda_2, \cdots, \lambda_n)\|^2)],$$

where

$$\mathcal{R}_1(\theta) = \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[1-\lambda_1]}{m} \sum_{i=1}^m (h'(f_{\theta}(x_i)) - y_i) \nabla f_{\theta}(x_i)^T \mathbb{E}_{r_x \sim \mathcal{D}_X}[r_x - x_i],$$

$$\mathcal{R}_{2}(\theta) = \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[\sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \nabla f_{\theta}(x_{i})^{T} \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})(r_{x} - x_{i})^{T}] \nabla_{\theta} f(x_{i})$$
$$+ \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[(1 - \lambda_{1})^{2} - \sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \nabla f_{\theta}(x_{i})^{T} \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})] \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})^{T}] \nabla_{\theta} f(x_{i})$$

$$\mathcal{R}_{3}(\theta) = \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[\sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})\nabla^{2}f_{\theta}(x_{i})(r_{x} - x_{i})^{T}] \\ + \frac{\mathbb{E}_{\boldsymbol{\lambda} \sim \tilde{\mathcal{D}}_{\Lambda}}[(1 - \lambda_{1})^{2} - \sum_{j=2}^{m} \lambda_{j}^{2}]}{2m} \sum_{i=1}^{m} (h'(f_{\theta}(x_{i})) - y_{i}) \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})]\nabla^{2}f_{\theta}(x_{i}) \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}}[(r_{x} - x_{i})^{T}].$$

# C Adversarial Robustness of MSDA

Let us scrutinize adversarial robustness in MSDA. We adopt the logistic loss, so  $l(\theta, z) = \log(1 + \exp(f_{\theta}(x))) - yf_{\theta}(x)$  where  $y \in \{0, 1\}$ . Define  $g(s) = e^s/(1 + e^s)$ . As [75], we scrutinize the logistic regression with  $f_{\theta}(x)$  as ReLU or leaky-ReLU network. Then, we have  $f_{\theta}(x) = \nabla f_{\theta}(x_i)^{\mathsf{T}} x_i$  and  $\nabla^2 f_{\theta}(x_i) = 0$ . We consider the adversarial loss with  $l_2$  attack of size  $\epsilon \sqrt{d}$  (d is the dimension of  $\theta$ ), that is,  $L_m^{\mathrm{adv}}(\theta) = \frac{1}{m} \sum_{i=1}^m \max_{\|\delta_i\|_2 \le \epsilon \sqrt{d}} l(\theta, (x_i + \delta_i, y_i))$ .

**Theorem 3.** In MSDA, we suppose that  $f_{\theta}(x) = \nabla f_{\theta}(x_i)^{\intercal} x_i, \nabla^2 f_{\theta}(x_i) = 0$  and there exists a constant  $c_x > 0$  that  $||x_i||_2 \le c_x \sqrt{d}$  for all  $i \in [m]$ . Then for any  $\theta \in \Theta$ , we have

$$\tilde{L}_m^{(MSDA)} \ge \frac{1}{m} \sum_{i=1}^m \tilde{l}_{adv}(\epsilon_i \sqrt{d}, z) \ge \frac{1}{m} \sum_{i=1}^m \tilde{l}_{adv}(\epsilon_{cut} \sqrt{d}, z),$$

where 
$$\epsilon_{cut} = c_x \min\left(\min_i |\cos(\nabla f_{\theta}(x_i), x_i)| \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}}[1-\lambda], \min_i \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}, M} s(M, x_i)\right)$$
 and  $s(M, x_i) = \frac{(\sqrt{1-M} \odot \nabla f(x_i))^T (\sqrt{1-M} \odot x_i)}{\|\nabla f(x_i)\|_2 \|x_i\|_2}.$ 

This bound appears to be fertile at first glance. However, as  $||f_{\theta}(x_i)||_2$  increases after training accuracy reaches 100%, the logistic loss decreases. Therefore, due to  $||f_{\theta}(x_i)||_2 = ||\nabla f_{\theta}(x_i)^{\mathsf{T}} x_i||_2 = ||\nabla f_{\theta}(x_i)||_2 ||x_i||_2 \cos(\nabla f_{\theta}(x_i), x_i), \cos(\nabla f_{\theta}(x_i), x_i)$  would be larger. Furthermore, under the CutMix case, since M distribution is relatively uniform,  $s(M, x_i)$  will be similar to  $\min_i |\cos(\nabla f_{\theta}(x_i), x_i)| \mathbb{E}_{\lambda \sim \tilde{D}_\lambda}[1 - \lambda]$  [75]. Moreover, empirical results support Mixup's or Cutmix's adversarial robustness [55, 70, 74].

proof of Theorem 3.

**Fact 1.** [75] The second order Taylor approximation of  $L_m^{adv}(\theta)$  is  $\frac{1}{m} \sum_{i=1}^m \tilde{l}_{adv}(\epsilon \sqrt{d}, z)$  where fore any  $\eta > 0, x \in \mathbb{R}^d$  and  $y \in \{0, 1\}$ ,

$$\tilde{l}_{adv}(\eta, z) = l(\theta, z) + \eta |g(f_{\theta}(x)) - y| \|\nabla f_{\theta}(x)\|_{2} + \frac{\eta^{2}d}{2} \cdot |h''(f_{\theta}(x))| \cdot \|\nabla f_{\theta}(x)\|_{2}^{2}.$$

We set  $\mathbb{E}_{r_x}(r_x) = 0$  by parallel translation. We define

$$s(M, x_i) = \frac{(\sqrt{1 - M} \odot \nabla f(x_i))^T (\sqrt{1 - M} \odot x_i)}{\|\nabla f(x_i)\|_2 \|x_i\|_2}.$$

For every MSDA, we have  $\mathcal{R}_1$  as

$$\mathcal{R}_1(\theta) = \frac{\mathbb{E}_{\lambda}(1-\lambda)}{m} \sum_{i=1}^m (y_i - g(f_{\theta}(x_i))) f_{\theta}(x_i),$$

and since  $\theta \in \Theta$ , we have  $(y_i - g(f_{\theta}(x_i)))f_{\theta}(x_i) \ge 0$ . Therefore,

$$\mathcal{R}_{1}(\theta) = \frac{\mathbb{E}_{\lambda}(1-\lambda)}{m} \sum_{i=1}^{m} |y_{i} - g(f_{\theta}(x_{i}))| |f_{\theta}(x_{i})|$$

$$= \frac{\mathbb{E}_{\lambda}(1-\lambda)}{m} \sum_{i=1}^{m} |y_{i} - g(f_{\theta}(x_{i}))| ||\nabla f_{\theta}(x_{i})||_{2} ||x_{i}||_{2} |\cos(\nabla f_{\theta}(x_{i}), x_{i})|$$

$$\geq \frac{\min_{i} ||x_{i}||_{2} \min_{i} |\cos(\nabla f_{\theta}(x_{i}), x_{i})| \mathbb{E}_{\lambda}(1-\lambda)}{m} \sum_{i=1}^{m} |y_{i} - g(f_{\theta}(x_{i}))| ||\nabla f_{\theta}(x_{i})||_{2}$$

Moreover, we can eliminate  $\mathcal{R}_3$  since  $\nabla^2 f_\theta = 0$ . So, we only focus on  $\mathcal{R}_2$  term. We have

$$\begin{aligned} \mathcal{R}_{2}^{(\text{MSDA})}(\theta) &= \frac{1}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbb{E}_{M}(1-M)^{\mathsf{T}} \mathbb{E}_{r_{x} \sim \mathcal{D}_{X}} \left( \nabla f(x_{i}) \odot (r_{x}-x_{i}) \left( \nabla f(x_{i}) \odot (r_{x}-x_{i}) \right)^{\mathsf{T}} \right) (1-M) \\ &\geq \frac{1}{2m} \sum_{i=1}^{m} h''(f_{\theta}(x_{i})) \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbb{E}_{M}(1-M)^{\mathsf{T}} \left( \left( \nabla f(x_{i}) \odot x_{i} \right) \left( \nabla f(x_{i}) \odot x_{i} \right)^{\mathsf{T}} \right) (1-M) \\ &= \frac{1}{2m} \sum_{i=1}^{m} |g(f_{\theta}(x_{i}))(1-g(f_{\theta}(x_{i})))| \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}} \mathbb{E}_{M}((\sqrt{1-M} \odot \nabla f(x_{i}))^{T} (\sqrt{1-M} \odot x_{i}))^{2} \\ &= \frac{1}{2m} \sum_{i=1}^{m} |g(f_{\theta}(x_{i}))(1-g(f_{\theta}(x_{i})))| \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}, M} s(M, x_{i})^{2} \| \nabla f(x_{i}) \|_{2}^{2} \|x_{i}\|_{2}^{2} \\ &\geq \frac{\min_{i} \|x_{i}\|^{2} \min_{i} \mathbb{E}_{\lambda \sim \tilde{D}_{\lambda}, M} s(M, x_{i})^{2}}{2m} \sum_{i=1}^{m} |g(f_{\theta}(x_{i}))(1-g(f_{\theta}(x_{i})))| \| \nabla f(x_{i}) \|_{2}^{2}, \end{aligned}$$

which concludes the theorem.

# **D** Generalization properties of MSDA

The data-dependent MSDA regularization can be altered to the original empirical risk minimization problem with a constrained function set. The Rademacher complexity of this constrained function set is  $\mathcal{O}(1/\sqrt{n})$ , which leads to the generalization properties of MSDA. We investigate two models. The first is the GLM model, which has the loss function  $l(\theta, z) = A(\theta^{\intercal}x) - y\theta^{\intercal}x$ . The second is two-layer ReLU networks, which can be parameterized as  $f_{\theta}(x) = \theta_1^{\intercal}\sigma(Wx) + \theta_0$ . In this case, we consider the mean square error (MSE) loss function  $(i.e., l(\theta) = \frac{1}{m} \sum_{i=1}^{m} (y_i - f_{\theta}(x_i))^2)$ 

#### D.1 GLM Model

For GLM, using (14), since the prediction of the GLM model is invariant to the scaling of the training data, we think the dataset  $\hat{D} = {\hat{z}_i}_{i=1}^m$  with  $\hat{x}_i = 1 \oslash \overline{M} \odot (M \odot x_i + (1 - M) \odot r_x)$  where  $\overline{M} = \mathbb{E}M$ . Then, the loss function is

$$L_m^{(\text{MSDA})} = \frac{1}{m} \mathbb{E}_{\lambda} \mathbb{E}_{r_x} \mathbb{E}_M \sum_{i=1}^m l(\theta, \tilde{z}_i) = \frac{1}{m} \mathbb{E}_{\xi} \sum_{i=1}^m (A(\hat{x}_i^{\mathsf{T}} \theta) - y_i \hat{x}_i^{\mathsf{T}} \theta),$$

where  $\xi$  denotes the randomness of  $\lambda, r_x$ , and M. By the second approximation of  $A(\cdot)$ , we can express  $A(\hat{x}_i^{\mathsf{T}}\theta)$  as

$$A(\hat{x}_i^{\mathsf{T}}\theta) = A(x_i^{\mathsf{T}}\theta) + A'(x_i^{\mathsf{T}}\theta)(\hat{x}_i - x_i)^{\mathsf{T}}\theta + \frac{1}{2}A''(x_i^{\mathsf{T}}\theta)\theta^{\mathsf{T}}(\hat{x}_i - x_i)(\hat{x}_i - x_i)^{\mathsf{T}}\theta$$

to approximate the loss function. Therefore, we have

$$\tilde{L}_{m}^{(\text{MSDA})} = \frac{1}{m} \sum_{i=1}^{m} A(x_{i}^{\mathsf{T}}\theta) + \frac{1}{m} \mathbb{E}_{\xi} \sum_{i=1}^{m} \left( A'(x_{i}^{\mathsf{T}}\theta)(\hat{x}_{i} - x_{i})^{\mathsf{T}}\theta + \frac{1}{2} A''(x_{i}^{\mathsf{T}}\theta)\theta^{\mathsf{T}}(\hat{x}_{i} - x_{i})(\hat{x}_{i} - x_{i})^{\mathsf{T}}\theta \right)$$
$$= \frac{1}{m} \sum_{i=1}^{m} A(x_{i}^{\mathsf{T}}\theta) + \frac{1}{m} \sum_{i=1}^{m} \left( \frac{1}{2} A''(x_{i}^{\mathsf{T}}\theta)\theta^{\mathsf{T}} \operatorname{Var}_{\xi}(\hat{x}_{i})\theta \right),$$
(26)

where  $\tilde{L}_m^{(\text{MSDA})}$  denotes the approximate loss of  $L_m^{(\text{MSDA})}$  since  $\mathbb{E}_{\xi} r_x = 0$  and  $\mathbb{E}_{\xi} \hat{x}_i = x_i$ . For calculating  $\text{Var}_{\xi}(\hat{x}_i)$ , we use the law of total variance. We have

$$\begin{aligned} \operatorname{Var}_{\xi}(\tilde{x_{i}}) &= \left(\frac{1}{\bar{M}}\frac{1}{\bar{M}}^{\mathsf{T}}\right) \odot \operatorname{Var}_{\xi}\left(M \odot x_{i} + (1-M) \odot r_{x}\right) \\ &= \left(\frac{1}{\bar{M}}\frac{1}{\bar{M}}^{\mathsf{T}}\right) \odot \left(\mathbb{E}(\operatorname{Var}\left(M \odot x_{i} + (1-M) \odot r_{x} \,|\, \lambda, M\right) + \operatorname{Var}\left(\mathbb{E}\left(M \odot x_{i} + (1-M) \odot r_{x} \,|\, \lambda, M\right)\right) \\ &= \left(\frac{1}{\bar{M}}\frac{1}{\bar{M}}^{\mathsf{T}}\right) \odot \left(\mathbb{E}(1-M)\hat{\Sigma}_{X}(1-M)^{\mathsf{T}} + x_{i}\operatorname{Var}(M)x_{i}^{\mathsf{T}}\right)\right) \\ &= \frac{1}{\lambda^{2}}\left(\mathbb{E}(1-M)\hat{\Sigma}_{X}(1-M)^{\mathsf{T}} + x_{i}\operatorname{Var}(M)x_{i}^{\mathsf{T}}\right)\right),\end{aligned}$$

where  $\hat{\Sigma}_X = \frac{1}{m} \sum_{i=1}^m x_i x_i^{\mathsf{T}}$  with some notational ambiguity that  $\frac{1}{M} := \vec{1} \oslash \bar{M}$ . In our setting  $\bar{M} = \bar{\lambda}\vec{1}$  where  $\bar{\lambda} = \mathbb{E}_{\lambda \sim \tilde{\mathcal{D}}_{\lambda}}[\lambda]$ . Now we think the related dual problem to the (26):

$$\mathcal{W}_{\gamma} = \left\{ x \to \theta^{\mathsf{T}} x, \text{ such that } \theta \text{ satisfying} \right.$$
$$\left( \mathbb{E}_{x} A''(\theta^{\mathsf{T}} x) \right) \cdot \left( \theta^{\mathsf{T}} \left( \mathbb{E} (1 - M) \Sigma_{X} (1 - M)^{\mathsf{T}} \right) \theta + \theta^{\mathsf{T}} \left( \left( x \operatorname{Var}(M) x^{\mathsf{T}} \right) \right) \theta \right) \le \gamma \right\}.$$

Here, we assume the  $(\mathbb{E}_x[A''(v^{\intercal}x)])^2 \ge \rho \mathbb{E}_x(v^{\intercal}x)^2$ , which is called  $\rho$ -retentiveness [1, 75].

**Theorem 4-(a)** (Restated). Define  $\Sigma_X^{(M)} = \mathbb{E}(1-M)\Sigma_X(1-M)^{\intercal}$ . Suppose  $A(\cdot)$  is  $L_A$  Lipschitz, and  $\mathcal{X}, \mathcal{Y}, \Theta$  are all bounded. There exist constants L, B > 0, such that for all  $\theta$  that  $\theta^{\intercal} x \in \mathcal{W}_{\gamma}$ , which is the regularization induced by MSDA, we have

$$L(\theta) \le L_m(\theta) + 2LL_A \frac{1}{\sqrt{n}} (\gamma/\rho)^{1/4} \left( \sqrt{tr\left(\left(\Sigma_X^{(M)}\right)^{\dagger} \Sigma_X\right) + rank(\Sigma_X)} \right) + B\sqrt{\frac{\log(1/\delta)}{2n}},$$

with probability at least  $1 - \delta$ .

*Proof.* Firstly, we calculate the empirical Rademacher complexity of  $W_{\gamma}$ . For *n* i.i.d. Rademacher random variables  $\xi_1, \ldots, \xi_n$ , the definition of the empirical Rademacher complexity gives

$$\operatorname{Rad}(\mathcal{W}_{\gamma}, n) = \mathbb{E}_{\xi_{i}} \sup_{(\mathbb{E}_{x}A^{\prime\prime}(\theta^{\intercal}x)) \cdot \left(\theta^{\intercal}\Sigma_{x}^{(M)}\theta + \theta^{\intercal}((x\operatorname{Var}(M)x^{\intercal}))\right)\theta \leq \gamma} \frac{1}{n} \sum_{i=1}^{n} \xi_{i}\theta^{\intercal}x_{i}$$
$$\leq \mathbb{E}_{\xi_{i}} \sup_{(\mathbb{E}_{x}A^{\prime\prime}(\theta^{\intercal}x)) \cdot \theta^{\intercal}\Sigma_{x}^{(M)}\theta \leq \gamma} \frac{1}{n} \sum_{i=1}^{n} \xi_{i}\theta^{\intercal}x_{i}.$$

Due to the  $\rho$ -retentiveness, we have

$$\operatorname{Rad}(\mathcal{W}_{\gamma}, n) \leq \mathbb{E}_{\xi_{i}} \sup_{(\theta^{\intercal} \Sigma_{X} \theta) \cdot \left(\theta^{\intercal} \Sigma_{X}^{(M)} \theta\right) \leq \gamma/\rho} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} \theta^{\intercal} x_{i}$$
$$\leq \mathbb{E}_{\xi_{i}} \left( \sup_{\theta^{\intercal} \Sigma_{X} \theta \leq \sqrt{\gamma/\rho}} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} \theta^{\intercal} x_{i} + \mathbb{E}_{\xi_{i}} \sup_{\theta^{\intercal} \Sigma_{X}^{(M)} \theta \leq \sqrt{\gamma/\rho}} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} \theta^{\intercal} x_{i} \right)$$

For the first part of the RHS, define  $\tilde{x}_i = \Sigma_X^{\dagger/2} x_i$  and  $v = \Sigma_X^{1/2} \theta$ . Then, we have

$$\mathbb{E}_{\xi_{i}} \sup_{\theta^{\intercal} \Sigma_{X} \theta \leq \sqrt{\gamma/\rho}} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} \theta^{\intercal} x_{i} = \mathbb{E}_{\xi_{i}} \sup_{\|v\|^{2} \leq \sqrt{\gamma/\rho}} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} v^{\intercal} \tilde{x}_{i}$$

$$\leq \frac{1}{n} (\gamma/\rho)^{1/4} \mathbb{E}_{\xi_{i}} \left\| \sum_{i=1}^{n} \xi_{i} \tilde{x}_{i} \right\| \leq \frac{1}{n} (\gamma/\rho)^{1/4} \sqrt{\mathbb{E}_{\xi_{i}}} \left\| \sum_{i=1}^{n} \xi_{i} \tilde{x}_{i} \right\|^{2}$$

$$= \frac{1}{n} (\gamma/\rho)^{1/4} \sqrt{\sum_{i=1}^{n} \tilde{x}_{i}^{\intercal} \tilde{x}_{i}}.$$

$$(-(M))^{\frac{1}{2}}$$

Similarly, by defining  $\check{x}_i = \left(\Sigma_X^{(M)}\right)^{\dagger/2} x_i$  and  $v = \left(\Sigma_X^{(M)}\right)^{1/2} \theta$ ,

$$\mathbb{E}_{\xi_{i}} \sup_{\theta^{\intercal} \Sigma_{X}^{(M)} \theta \leq \sqrt{\gamma/\rho}} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} \theta^{\intercal} x_{i} = \mathbb{E}_{\xi_{i}} \sup_{\|v\|^{2} \leq \sqrt{\gamma/\rho}} \frac{1}{n} \sum_{i=1}^{n} \xi_{i} v^{\intercal} \check{x}_{i}$$
$$\leq \frac{1}{n} (\gamma/\rho)^{1/4} \mathbb{E}_{\xi_{i}} \left\| \sum_{i=1}^{n} \xi_{i} \check{x}_{i} \right\| \leq \frac{1}{n} (\gamma/\rho)^{1/4} \sqrt{\mathbb{E}_{\xi_{i}}} \left\| \sum_{i=1}^{n} \xi_{i} \check{x}_{i} \right\|^{2}$$
$$= \frac{1}{n} (\gamma/\rho)^{1/4} \sqrt{\sum_{i=1}^{n} \check{x}_{i}^{\intercal} \check{x}_{i}}.$$

Therefore,

$$\begin{aligned} \operatorname{Rad}(\mathcal{W}_{\gamma}) &= \mathbb{E}[\operatorname{Rad}(\mathcal{W}_{\gamma}, n)] \leq \frac{1}{n} (\gamma/\rho)^{1/4} \left( \sqrt{\sum_{i=1}^{n} \mathbb{E}_{x} \tilde{x}_{i}^{\mathsf{T}} \tilde{x}_{i}} + \sqrt{\sum_{i=1}^{n} \mathbb{E}_{x} \tilde{x}_{i}^{\mathsf{T}} \tilde{x}_{i}} \right) \\ &\leq \frac{1}{\sqrt{n}} (\gamma/\rho)^{1/4} \left( \sqrt{\operatorname{tr}\left( \left( \Sigma_{X}^{(M)} \right)^{\dagger} \Sigma_{X} \right)} + \operatorname{rank}(\Sigma_{X}) \right). \end{aligned}$$

The relationship between Rademacher complexity and generalization error [2] indicates Theorem 4.  $\Box$ 

#### D.2 Two-layer ReLU Networks

We perform MSDA on the final layer of the two-layer ReLU networks. Therefore, the setting is the same as Appendix D.1 with covariates  $\sigma(w_j^T x)$ . Due to the scaling of  $\theta_1$  and  $\theta_0$ , we consider training

 $\theta_1$  and  $\theta_0$  on the covariates  $1 \oslash \overline{M} \odot (M \odot (\sigma(Wx_i) - \overline{\sigma}_W) + (1 - M) \odot (\sigma(Wr_x) - \overline{\sigma}_W))$  where  $\overline{\sigma}_W = \frac{1}{n} \sum_{i=1}^m \sigma(Wx_i)$ . Putting GLM loss with  $A(\cdot) = \frac{1}{2} (\cdot)^2$ , we have

$$\tilde{L}_{m}^{(\text{MSDA})} = \frac{1}{m} \sum_{i=1}^{m} (f_{\theta}(x_{i}) - y_{i})^{2} + \frac{1}{m} \sum_{i=1}^{m} \left( \frac{1}{2} \theta^{\mathsf{T}} \operatorname{Var}_{\xi}(\sigma(W(M \odot x_{i} + (1 - M) \odot r_{x}))\theta) \right)$$
(27)

where  $\tilde{L}_m^{\rm (MSDA)}$  denotes the approximate loss of  $L_m^{\rm (MSDA)}$  and

$$\operatorname{Var}_{\xi}(\tilde{x}_{i}) = \frac{1}{\bar{\lambda}^{2}} \left( \mathbb{E}(1-M)\hat{\Sigma}_{X}^{\sigma}(1-M)^{\mathsf{T}} + \sigma(Wx_{i})\operatorname{Var}(M)\sigma(Wx_{i})^{\mathsf{T}}) \right)$$

where  $\hat{\Sigma}_X^{\sigma} = \operatorname{Var}_{r_x \sim \mathcal{D}_X} \sigma(Wr_x)$  with some notational ambiguity that  $\frac{1}{M} := \vec{1} \oslash \bar{M}$ . Now we think of the related dual problem to the equation 27:

$$\mathcal{W}_{\gamma}^{\mathrm{NN}} = \left\{ x \to f_{\theta}(x) = \theta_{1}^{\mathsf{T}} \sigma(Wx) + \theta_{0}, \text{ such that } \theta \text{ satisfying} \\ \theta_{1}^{\mathsf{T}} \left( \mathbb{E}(1-M) \Sigma_{X}^{\sigma} (1-M)^{\mathsf{T}} \right) \theta_{1} + \theta_{1}^{\mathsf{T}} \left( \mathbb{E}_{x} \left( \sigma(Wx) \operatorname{Var}(M) \sigma(Wx)^{\mathsf{T}} \right) \right) \theta_{1} \leq \gamma \right\},$$

where  $\Sigma_X^{\sigma} = \operatorname{Var}_x \sigma(Wx)$ .

**Theorem 4-(b)** (Restated). Define  $\Sigma_X^{\sigma,(M)} = \mathbb{E}(1-M)\Sigma_X^{\sigma}(1-M)^{\mathsf{T}}$ . Suppose  $\mathcal{X}, \mathcal{Y}, \Theta$  are all bounded. There exists constants L, B > 0, such that for all  $\theta$  that  $f_{\theta}(x) \in W_{\gamma}^{NN}$ , which is the regularization induced by Manifold MSDA, we have

$$L(\theta) \le L_m(\theta) + 4L\sqrt{\frac{\gamma\left(\operatorname{rank}(\Sigma_X^{\sigma,(M)}) + \left\|\left(\Sigma_X^{\sigma,(M)}\right)^{\dagger/2}\mu_{\sigma}\right\|^2\right)}{n}} + B\sqrt{\frac{\log(1/\delta)}{2n}},$$

with probability at least  $1 - \delta$ .

*Proof.* Firstly, we calculate the empirical Rademacher complexity of  $W_{\gamma}^{NN}$ . For the *n* i.i.d. Rademacher random variables  $\xi_1, \ldots, \xi_n$ , the definition of the empirical Rademacher complexity gives

$$\operatorname{Rad}(\mathcal{W}_{\gamma}^{\operatorname{NN}},n) = \mathbb{E}_{\xi_{i}} \sup_{\substack{\theta_{1}^{\mathsf{T}}\left(\mathbb{E}(1-M)\Sigma_{X}^{\sigma}(1-M)^{\mathsf{T}}\right)\theta_{1}+\theta_{1}^{\mathsf{T}}\left(\mathbb{E}_{x}(\sigma(Wx)\operatorname{Var}(M)\sigma(Wx)^{\mathsf{T}})\right)\theta_{1}\leq\gamma} \frac{1}{n} \sum_{i=1}^{n} \xi_{i}\theta_{1}^{\mathsf{T}}\sigma(Wx_{i})$$

$$\leq \mathbb{E}_{\xi_{i}} \sup_{\substack{\theta_{1}^{\mathsf{T}}\left(\mathbb{E}(1-M)\Sigma_{X}^{\sigma}(1-M)^{\mathsf{T}}\right)\theta_{1}\leq\gamma} \frac{1}{n} \sum_{i=1}^{n} \xi_{i}\theta_{1}^{\mathsf{T}}\sigma(Wx_{i})$$

$$\leq \mathbb{E}_{\xi_{i}} \sup_{\substack{\theta_{1}^{\mathsf{T}}\left(\mathbb{E}(1-M)\Sigma_{X}^{\sigma}(1-M)^{\mathsf{T}}\right)\theta_{1}\leq\gamma} \frac{1}{n} \sum_{i=1}^{n} \xi_{i}\theta_{1}^{\mathsf{T}}(\sigma(Wx_{i})-\mu_{\sigma}) + \mathbb{E}_{\xi_{i}} \sup_{\substack{\theta_{1}^{\mathsf{T}}\left(\mathbb{E}(1-M)\Sigma_{X}^{\sigma}(1-M)^{\mathsf{T}}\right)\theta_{1}\leq\gamma} \frac{1}{n} \sum_{i=1}^{n} \xi_{i}\theta_{1}^{\mathsf{T}}\mu_{\sigma},$$

where  $\mu_{\sigma} = \mathbb{E}[\sigma(Wx)]$ . Setting  $\tilde{\theta}_1^{\mathsf{T}} = \left(\Sigma_X^{\sigma,(M)}\right)^{1/2}$ , same technique with the proof of Theorem 3-(a) gives

$$\operatorname{Rad}(\mathcal{W}_{\gamma}^{\operatorname{NN}}) \leq 2\sqrt{\frac{\gamma\left(\operatorname{rank}(\boldsymbol{\Sigma}_{X}^{\sigma,(M)}) + \left\|\left(\boldsymbol{\Sigma}_{X}^{\sigma,(M)}\right)^{\dagger/2} \mu_{\sigma}\right\|^{2}\right)}{n}}$$

Finally, the relationship between Rademacher complexity and generalization error [2] indicates Theorem 4.  $\hfill \Box$ 

# E Extending Chidambaram *et al.* [12]

Chidambaram *et al.* [12] gave the theoretical analysis of Mixup. We follow this paper by using the unified framework that we used. By modifying several definitions, we get similar results with Chidambaram *et al.* [12].

Here, we assume that k classes have disjoint support, *i.e.*  $X = \bigcup_{i=1}^{k} X_i$  and  $X_i$  are mutually disjoint for i = 1, ..., k, where  $X_i$  is the support of the *i*th class. We consider cross-entropy loss. We define an associated probability measure  $\mathbb{P}_X$ . Then,  $L^{\text{MSDA}}$ , which is the expected loss for MSDA, can be expressed with

$$L^{(\text{MSDA})}(\theta) = \mathbb{E}_{z_1, z_2 \sim X} \mathbb{E}_{\lambda \sim \mathcal{D}_{\lambda}} \mathbb{E}_M l(\theta, \tilde{z}_{z_1, z_2}^{(\text{MSDA})}(\lambda, 1 - \lambda)),$$

where l is the cross entropy function. We can express  $L^{(\text{MSDA})}(\theta) = \sum_{i=1}^{k} \sum_{j=1}^{k} L_{i,j}^{(\text{MSDA})}(\theta)$  with  $i, j \in [k]$ , where  $L_{i,j}^{(\text{MSDA})}(\theta)$  is defined as

$$L_{i,j}^{(\mathrm{MSDA})}(\theta) = \mathbb{E}_{z_1, z_2 \sim X} \mathbb{E}_{\lambda \sim \mathcal{D}_{\lambda}} \mathbb{E}_M \left[ l(\theta, \tilde{z}_{z_1, z_2}^{(\mathrm{MSDA})}(\lambda, 1 - \lambda)) I(z_1 \in X_i, z_2 \in X_j) \right].$$

 $L_{i,j}^{(MSDA)}(\theta)$  is the full MSDA cross entropy loss corresponding to the mixing points from classes i and j. The goal of standard training is to learn a classifier  $h \in \arg \min_{g \in \mathcal{C}} L(g, \mathbb{P}_X)$  where  $\mathcal{C}$  is the classifier set. Any such classifier h will satisfy  $h(x)_i = 1$  on  $X_i$  since the  $X_i$  are disjoint.

We modify some definitions in [12, Section 2.2]. We define  $A_{x,\epsilon}^{i,j}$  and  $A_{x,\epsilon,\delta}^{i,j}$  as

$$\begin{split} A_{x,\epsilon}^{i,j} &= \{(s,t,\lambda,M) \in X_i \times X_j \times [0,1] \times \mathbb{R}^n : M \odot s + (1-M) \odot t \in B_{\epsilon}(x)\} \\ A_{x,\epsilon,\delta}^{i,j} &= \{(s,t,\lambda,M) \in X_i \times X_j \times [0,1-\delta] \times \mathbb{R}^n : M \odot s + (1-M) \odot t \in B_{\epsilon}(x)\} \\ X_{\text{MSDA}} &= \left\{ x \in \mathbb{R}^n : \bigcup_{i,j} A_{x,\epsilon}^{i,j} \text{ has positive measure for every } \epsilon > 0 \right\} \\ \xi_{x,\epsilon}^{i,j} &= \mathbb{E}_{z_1, z_2 \sim X} \mathbb{E}_{\lambda \sim \mathcal{D}_{\lambda}} \mathbb{E}_M [I(z_1 \in X_i, z_2 \in X_j)] \\ \xi_{x,\epsilon,\lambda}^{i,j} &= \mathbb{E}_{z_1, z_2 \sim X} \mathbb{E}_{\lambda \sim \mathcal{D}_{\lambda}} \mathbb{E}_M [\lambda I(z_1 \in X_i, z_2 \in X_j)]. \end{split}$$

**Definition E.1** ([12]). Let  $C^*$  to be the subset of C for which every  $h \in C^*$  satisfies  $h(x) = \lim_{\epsilon \to 0} \arg\min_{\theta \in [0,1]} L^{(\text{MSDA})}(\theta)|_{B_{\epsilon}(x)}$  for all  $x \in X_{\text{MSDA}}$  when the limit exists. Here,  $L^{(\text{MSDA})}(\theta)|_{B_{\epsilon}(x)}$  represents the MSDA loss for a constant function with value  $\theta$  with the restriction of each term in  $L^{(\text{MSDA})}$  to the set  $A_{x,\epsilon}^{i,j}$ 

 $C^*$  includes deep neural networks. Below lemmas and theorems can be proved with the same technique as Chidambaram *et al.*.

**Lemma 1.** Any function  $h \in \arg \min_{g \in C^*} L^{(MSDA)}(g, \mathbb{P}_X, \mathbb{P}_f)$  satisfies  $L^{(MSDA)}(h) \leq L^{(MSDA)}(g)$  for any continuous  $g \in C$ 

**Theorem 6.** For any point  $x \in X_{MSDA}$  and  $\epsilon > 0$ , there exists a continuous function  $h_{\epsilon}$  satisfying

$$h_{\epsilon}^{i}(x) = \frac{\xi_{x,\epsilon}^{i,i} + (\sum_{j \neq i} \xi_{x,\epsilon,\lambda}^{i,j} + (\xi_{x,\epsilon}^{j,i} - \xi_{x,\epsilon,\lambda}^{j,i}))}{\sum_{q=1}^{k} \xi_{x,\epsilon}^{q,q} + \sum_{j \neq q} (\xi_{x,\epsilon,\lambda}^{q,j} + (\xi_{x,\epsilon}^{j,q} - \xi_{x,\epsilon,\lambda}^{j,q}))}$$

and its limits exist when  $\epsilon \to 0$ .

We give an assumption: for a point  $x \in X_{MSDA}$ , there exists a class *i* that *x* is closest to  $X_i$  for arbitrary MSDA expression of *x*, and *x* cannot be expressed by MSDA expression between non-*i* classes. The formal assumption and its geometric intuition can be found in [12].

**Theorem 7.** If x satisfies the above assumption, with respect to a class i, then for every  $h \in \arg \min_{a \in C^*} L_{MSDA}(g)$ , we have that h classifies x as the class i and h is continuous at x.

Theorem 7 indicates that if we observe a new sample that can satisfy the above assumption with the class *i*, which is also a distance up to  $\min_j d(X_i, X_j)/2$  from the class *i*, the model trained with MSDA will classify it as *i*. Therefore, Theorem 7 is closely related to the generalization properties of MSDA. All proof of this section is identical to Chidambaram *et al.*.



Figure 6: Comparison of regularized partial gradients between Mixup and CutMix along the sorted pixel indices.

 Table 4: HMix performances in both scenarios of Table 1.

	Mixup	CutMix	$\Delta \left( {\rm CutMix} \text{ - Mixup} \right)$	HMix ( <i>r</i> =0.5)	HMix ( <i>r</i> =0.75)
Scenario 1: Large crop	58.3	64.4	+6.1	61.3 (+3.0 vs. Mixup)	<b>63.7</b> (+ <b>5.4</b> vs. Mixup)
Scenario 2: Small crop	67.7	67.0	-0.7	67.6 (+0.6 vs. CutMix)	<b>67.2</b> (+ <b>0.2</b> vs. CutMix)

### **F** Experimental Details

**Regularized input gradients experiments (Figure 4)** We investigate the amount of the regularized input gradients by  $|\partial_v f_\theta(x) \partial_{v+p} f_\theta(x)|$  with respect to the pixel distance vector p. We then compute the partial gradients for the validation images x, which have not been seen during training, and normalize them to sum to 1 as

 $\operatorname{PartialGradProd}(x,p) = \max_{v} |\partial_{v} f_{\theta}(x) \partial_{v+p} f_{\theta}(x)|$ 

for different  $f_{\theta}$  trained by different MSDA methods. Finally, we visualize the pixel-wise maximum values of PartialGradProd(x, p) for the validation images x in Figure 4. We train ResNet-50 [26] models on the ImageNet-1K dataset with  $64 \times 64$  image size. We show additional visualization in Figure 6, where the x-axis denotes the pixel indices sorted by the pixel distance ||p||. As shown in Figure 4 and Figure 6, CutMix more regularizes the partial gradient products  $|\partial_v f_{\theta}(x) \partial_{v+p} f_{\theta}(x)|$  for a closer p than Mixup.

**CIFAR-100 experiments** We follow the experimental setting of [39] and utilize [39]'s codebase<sup>2</sup>. We train all networks for 300 epochs using the SGD optimizer with a learning rate of 0.2 and a batch size of 100. The learning rate is decayed by a factor of 0.1 at 100 and 200 epochs. We set the hyper-parameter  $\alpha$  for Mixup, CutMix, and Stochastic Mixup & CutMix to 1.  $\alpha$  for HMix and GMix were set to 1 and 0.5, respectively. We use r = 0.5 for HMix. All the experiments in CIFAR-100 were repeated three times, and we report the average accuracy.

**ImageNet-1K experiments** We utilize *timm*<sup>3</sup> pytorch codebase. We train ResNet-50 [26] for 300 epochs using the SGD optimizer with a learning rate of 0.1, weight decay of  $2 \times 10^{-5}$ , and batch size of 512. We use cosine learning rate scheduling. We set the hyper-parameter  $\alpha$  for all methods except Mixup to 1, while Mixup has  $\alpha = 0.8$ . We use r = 0.75 for HMix.

# **G** Additional Experiments

**HMix in the scenarios of Table 1.** We conduct HMix on two scenarios in Table 1. Results are in Table 4. We use r = 0.5 and r = 0.75 for HMix. HMix shows consistently better performance than Mixup and CutMix for the scenarios 1 and 2, respectively, with meaningful performance gaps. Through the results, we can empirically confirm our proposed method enjoyed the advantages of Mixup and CutMix.

<sup>&</sup>lt;sup>2</sup>https://github.com/snu-mllab/PuzzleMix

<sup>&</sup>lt;sup>3</sup>https://github.com/rwightman/pytorch-image-models

Augmentation Method	ImageNet-1K	ImageNet-occ	ImageNet-C	FGSM
Vanilla (no MDSA)	75.68 (+0.00)	55.26 (+0.00)	42.57 (+0.00)	8.55 (+0.00)
Mixup	77.78 (+2.10)	60.34 (+5.08)	51.73 (+9.16)	27.78 (+19.23)
CutMix	78.04 (+2.36)	71.51 (+16.25)	44.18 (+1.55)	33.63 (+25.08)
HMix (ours)	78.38 (+2.70)	71.13 (+15.87)	46.37 (+3.80)	34.98 (+26.44)
GMix (ours)	78.13 (+2.45)	62.76 (+7.50)	45.97 (+3.40)	31.02 (+21.47)

Table 5: Robustness benchmarks. Comparison of various MSDA methods on ResNet-50 architecture.

Table 6: Ablation study on hyper-parameters.

(a) Impact on a for CMin

(a) Impact on r for HMix.		(b) Impact on $\alpha$ for HMix.			
PreActRN18	CIFAR-100 acc	PreActRN18	CIFAR-100 acc	PreActRN18	CIFAR-100 acc
Vanilla (no MDSA) Mixup ( $r = 0$ ) CutMix ( $r = 1.0$ ) <b>HMix</b> ( $r = 0.75$ ) <b>HMix</b> ( $r = 0.5$ )	76.73 77.21 78.66 <b>79.43</b> 79.25	Vanilla (no MDSA) Mixup ( $\alpha = 1.0$ ) CutMix ( $\alpha = 1.0$ ) <b>HMix</b> ( $\alpha = 0.5$ ) <b>HMix</b> ( $\alpha = 1.0$ )	76.73 77.21 78.66 78.47 <b>79.25</b>	Vanilla (no MDSA) Mixup ( $\alpha = 1.0$ ) CutMix ( $\alpha = 1.0$ ) <b>GMix</b> ( $\alpha = 0.25$ ) <b>GMix</b> ( $\alpha = 0.5$ )	76.73 77.21 78.66 78.60 <b>79.17</b>
<b>HMix</b> $(r = 0.0)$	78.05	<b>HMix</b> ( $\alpha = 2.0$ )	78.85	<b>GMix</b> ( $\alpha = 0.75$ ) <b>GMix</b> ( $\alpha = 1.0$ )	78.64 79.05

**Robustness benchmarks.** As observed by the previous study [13], the different choice of MSDA methods also affects the extreme case of the test samples, *e.g.*, distribution shifts. In this experiments, we provide an understanding of the relationship between MSDA and various test scenarios and through the lens of our theoretical analysis. We conduct various MSDA methods on robustness benchmarks such as ImageNet-1K occlusion accuracy (center occluded images following [13, 70]), ImageNet-C accuracy [28] and adversarially attacked ImageNet test accuracy by FGSM attack [22]. Results are in Table 5. We use the same experimental settings of the ImageNet classification in Table 3. Overall results show that HMix and GMix are located in between CutMix and Mixup.

CutMix performs better than Mixup in the occlusion accuracy. Here, the local occluded areas have no information to distinguish objects, but other local areas are informative. Hence, it is important to capture shorter-relationship rather than global-relationship. Thus we can expect that CutMix is better than Mixup, and not surprisingly, HMix and GMix are located in between CutMix and Mixup.

ImageNet-C style corruptions (*e.g.*, adding Gaussian noise for the entire image) distort the local information. In this case, the shorter-distance relationships are significantly damaged and sometimes useless to distinguish the object, hence the longer-distance relationships are more important. Hence, we can expect Mixup works better than CutMix in ImageNet-C. As HMix and GMix less weight shorter-distance relationships than CutMix (but more weight than Mixup), we can observe that HMix and GMix ImageNet-C performances are better than CutMix, but worse than Mixup.

Ablation study on hyper-parameters. We study the impacts of hyper-parameters  $\alpha$  and r for HMix and GMix on CIFAR-100 classification. We use the PreActResNet-18 with the same experimental setting as in Table 2. Results are in Table 6. We highlight the results with our hyper-parameter choices (r = 0.5,  $\alpha = 1.0$  for HMix,  $\alpha = 0.5$  for GMix) in the gray cells. For r, we find that HMix with r = 0.75 performs better than r = 0.5 with a marginal gap (+0.18%p). For  $\alpha$ , our hyper-parameters show the best performance. Overall results confirm that HMix and GMix are not sensitive to those hyper-parameters and consistently show better or compatible performance against Mixup and CutMix.

# **H** Regularizer coefficients of HMix

Define

$$h(x,s) = \min(x, n-s), \qquad l(x,s) = \max(x-s,0),$$

$$a_{jk,s} = \frac{\max(\min(h(j_1,s) - l(k_1,s), h(k_1,s) - l(j_1,s)), 0) \max(\min(h(j_2,s) - l(k_2,s), h(k_2,s) - l(j_2,s)), 0)}{(n-s)^2},$$

$$(28)$$

$$o_s = \frac{\lambda n^2}{n^2 - s^2},$$

$$v(p,s) = \frac{(h(p_1,s) - l(p_1,s)) * (h(p_2,s) - l(p_2,s))}{(n-s)^2}.$$

Note that (28) is extension of (10) by putting  $s = [\sqrt{1 - \lambda}n]$  in (28). Then, HMix with hyperparameter r has regularizer coefficient  $a_{ij}$  as

$$s = [\sqrt{1 - \lambda}\sqrt{rn}]$$
  
$$a_{ij} = o(s)(1 - o(s))(v(i, s) + v(j, s)) + o(s)a_{ij,s} + (1 - o(s))(1 - o(s)).$$

We plotted this value in Figure 4 when r = 0.7.

### I What MSDA can be applied in our Thereoms?

Our theorems can be applied to any MSDA method with an analogous formula, regardless of the assumption of the shape of the mask. In this paper, we mainly focused on Mixup and CutMix because they are the most common MSDA methods among the whole MSDA family as well as their behaviors are distinctly different in terms of our theorem. In this section, we note several nontrivial remarks for understanding the setup of our paper.

**ResizeMix.** ResizeMix [54] can be explained by our theorem if we add the assumption on the dataset  $\mathcal{D}_X$  that  $\mathcal{D}_X$  has all resized versions of the image. ResizeMix uses the resized version of input (i.e., one of the mixed patches is the "resized" version, not a cropped one) where the random resize is applied to the whole dataset. In other words, ResizeMix is a special case of CutMix when we apply a special version of random resize crop operation. Hence, if we assume a different version of random resize crop rather than the standard version (independent of our theoretical results and underlying assumptions), ResizeMix is equivalent to CutMix, which leads to the same theoretical result as CutMix.

**FMix.** FMix [25] randomly samples the mask from the Fourier space. Since FMix is one of the static MSDA, we can directly apply our Theorem 1-4.

**SaliencyMix, PuzzleMix, and Co-Mixup.** SaliencyMix [64] uses saliency map to generate new MSDA sample. PuzzleMix [39], and Co-Mixup [38] are dynamic MSDA, where they use saliency map and transport. These methods give a state-of-the-art performance. Theorem 1 can deal with this problem, but hard to interpret. To be specific, the second equality of Equation (8) do not hold anymore; it is hard to interpret the approximated loss function as an input gradient / Hessian regularizer.

**StyleMix and Manifold Mixup.** StyleMix [29] uses pre-trained style encoder and decoder. StyleMix linearly mixes content and style. Manifold Mixup [65] mixes samples in the feature level. Therefore, the theorems cannot be directly applied to StyleMix and Manifold Mixup.

AutoMix. Recently, AutoMix [45] gives state-of-the-art results. AutoMix utilized joint loss to generate the mask M: classification loss and generation loss for training M. Therefore, the mask depends on the mixing samples, indicating that AutoMix is a dynamic MSDA. As in previous paragraphs, Theorem 1 holds, but it is not easy to interpret each term.