# Learning with Multiplicative Perturbations

Xiulong Yang, and Shihao Ji
Department of Computer Science
Georgia State University, USA
xyang22@gsu.edu; sji@gsu.edu

*Abstract*—Adversarial Training (AT) and Virtual Adversarial Training (VAT) are the regularization techniques that train Deep Neural Networks (DNNs) with adversarial examples generated by adding small but worst-case perturbations to input examples. In this paper, we propose xAT and xVAT, new adversarial training algorithms that generate multiplicative perturbations to input examples for robust training of DNNs. Such perturbations are much more perceptible and interpretable than their additive counterparts exploited by AT and VAT. Furthermore, the multiplicative perturbations can be generated transductively or inductively, while the standard AT and VAT only support a transductive implementation. We conduct a series of experiments that analyze the behavior of the multiplicative perturbations and demonstrate that xAT and xVAT match or outperform state-of-the-art classification accuracies across multiple established benchmarks while being about 30% faster than their additive counterparts. Our source code can be found at https://github.com/sndnyang/xvat

## I. INTRODUCTION

Over the past few years, Deep Neural Networks (DNNs) have achieved state-of-the-art performance on a wide range of learning tasks. However, the success of DNNs has a high reliance on large sets of labeled examples; when trained on small datasets, DNNs plague to overfitting (if not regularized properly). For many practical applications, collecting a large amount of labeled examples is very expensive and/or time-consuming. To address this issue, researchers have investigated a host of techniques, such as Dropout [1], AT [2, 3], VAT [4], and Mixup [5], to regularize the training of DNNs. Such techniques usually augment the loss function of DNNs with a regularization term to prevent the model from overfitting when the labeled train set is small.

Several studies have found that the performance of DNNs can be improved significantly by enforcing the prediction consistency of DNNs in response to original inputs and their perturbated versions. For instance, Szegedy et al. [2] have found that very tiny perturbations to input samples (a.k.a., adversarial examples) can easily fool a well-trained DNN because the decision boundary of the DNN can change sharply around some data points. To improve the robustness of DNNs, they introduce AT to regularize the training of DNNs by augmenting the training set with adversarial examples. Furthermore, VAT [4, 6] extends the adversarial training principle of AT from supervised learning to semi-supervised learning by generating adversarial perturbations on unlabeled examples based on a divergence measure. However, the perturbations exploited by AT and VAT are *additive* in the sense that these perturbations are added pixel-wise to input examples.

In this paper, we propose a new type of perturbations called *multiplicative* perturbations that are generated via an $L_0$-norm regularized optimization and are *multiplied* to input examples pixel by pixel. This is in a stark contrast to the additive perturbations exploited by AT and VAT as the additive perturbations are generated by maximizing a divergence measure and *added* to input examples pixel by pixel. To illustrate the differences, Fig. 1 demonstrates the learning pipelines of the additive perturbations and our multiplicative perturbations, with the main differences highlighted in the dashed boxes, where a pair of forward and backward propagations in the additive pipeline is replaced by a sparse mask generator in the multiplicative pipeline. Given an input image, the sparse mask generator outputs an adversarial mask, which is subsequently multiplied to the original input to generate a multiplicative adversarial example. We optimize the sparse mask generator *adversarially* to maximize a divergence measure under an $L_0$-norm regularization. Similar to the additive perturbations, the multiplicative perturbations can be generated for labeled examples and unlabeled examples, and therefore can be used for supervised learning and semi-supervised learning. In light of the similarity to AT and VAT, we call our multiplicative AT and VAT as xAT and xVAT, with x denoting multiplication. Furthermore, our method can generate multiplicative perturbations transductively or inductively, while the additive perturbations are optimized transductively through backpropagation. For the reason to be discussed later, the parameters of sparse mask generator and the classifier of xAT and xVAT can be optimized simultaneously in one step, while AT and VAT have to optimize the additive perturbations and the classifier alternatively in two steps. As a result, xAT and xVAT are computationally more efficient than their additive counterparts.

Our main contributions are summarized as follows:
1) We introduce a new type of perturbations for robust training of DNNs that are multiplicative instead of additive; compared to the conventional additive perturbations, the multiplicative ones are much more perceptible and interpretable;
2) The multiplicative perturbations can be generated transductively or inductively, and the sparse mask generator and the classifier can be optimized simultaneously in one step, making xAT and xVAT computationally more efficient than their additive counterparts;
3) On four image classification benchmarks, xAT and xVAT match or outperform the state-of-the-art algorithms while being about 30% faster.
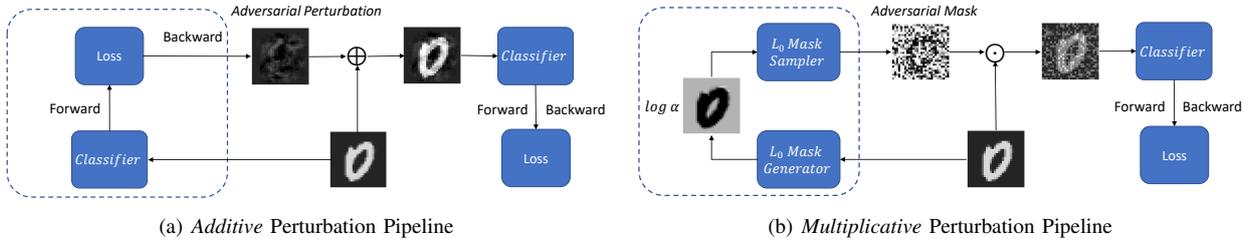
(a) *Additive* Perturbation Pipeline       (b) *Multiplicative* Perturbation Pipeline

Fig. 1. Comparison of the *additive* perturbation pipeline and the *multiplicative* perturbation pipeline.

## II. METHOD

Assume that we have a labeled dataset $\mathcal{D}_l = \{(\boldsymbol{x}_l^i, y_l^i),\ i = 1, 2, \cdots, N_l\}$, and an unlabeled dataset $\mathcal{D}_u = \{\boldsymbol{x}_u^j,\ j = 1, 2, \cdots, N_u\}$, where $\boldsymbol{x}^k \in \mathbb{R}^P$ denotes the $k$-th input sample and $y^k$ is its corresponding label. We use $p(y|\boldsymbol{x}, \boldsymbol{\theta})$ to denote the output distribution of a classifier, parameterized by $\boldsymbol{\theta}$, in response to an input $\boldsymbol{x}$. In supervised learning, we optimize model parameter $\boldsymbol{\theta}$ by minimizing an empirical risk on $D_l$, while in semi-supervised learning both $D_l$ and $D_u$ are utilized to optimize model parameter $\boldsymbol{\theta}$.

### A. Additive Perturbations with AT and VAT

AT [3] and VAT [4] are two regularization techniques that have been proposed to generate small but worst-case perturbations for robust training of DNNs. Specifically, AT [3] solves the following constrained optimization problem:

$$\mathcal{L}_{\text{AT}}(\boldsymbol{x}_l, y_l, \boldsymbol{r}_{\text{adv}}, \boldsymbol{\theta}) = D\left[h(y_l|\boldsymbol{x}_l), p(y|\boldsymbol{x}_l + \boldsymbol{r}_{\text{adv}}, \boldsymbol{\theta})\right] \quad (1)$$
$$\text{with } \boldsymbol{r}_{\text{adv}} = \underset{\boldsymbol{r};\|\boldsymbol{r}\| \leq \epsilon}{\arg\max} D\left[h(y_l|\boldsymbol{x}_l), p(y|\boldsymbol{x}_l + \boldsymbol{r}, \boldsymbol{\theta})\right],$$

where $D[p, q]$ is a divergence measure between two distributions $p$ and $q$. For the task of image classification, $p$ and $q$ are the probability vectors whose $i$-th element denotes the probability of an input image belonging to class $i$. In AT, $D$ is the cross entropy loss $D[p, q] = -\sum_i p_i \log q_i$ and $h(y|\boldsymbol{x})$ is the one-hot encoding of label $y$ for sample $\boldsymbol{x}$. Since $h(y|\boldsymbol{x})$ requires the true label $y$ of $\boldsymbol{x}$, AT can only be applied to supervised learning. To extend the adversarial training to unlabeled samples, VAT [4] substitutes $h(y|\boldsymbol{x})$ with the predicted classification probability $p(y|\boldsymbol{x}, \boldsymbol{\theta})$ and solves a slightly different constrained optimization problem:

$$\mathcal{L}_{\text{VAT}}(\boldsymbol{x}_*, \boldsymbol{r}_{\text{adv}}, \boldsymbol{\theta}) = D\left[p(y|\boldsymbol{x}_*, \boldsymbol{\theta}), p(y|\boldsymbol{x}_* + \boldsymbol{r}_{\text{adv}}, \boldsymbol{\theta})\right] \quad (2)$$
$$\text{with } \boldsymbol{r}_{\text{adv}} = \underset{\boldsymbol{r};\|\boldsymbol{r}\|_2 \leq \epsilon}{\arg\max} D\left[p(y|\boldsymbol{x}_*, \boldsymbol{\theta}), p(y|\boldsymbol{x}_* + \boldsymbol{r}, \boldsymbol{\theta})\right],$$

where $\boldsymbol{x}_*$ can be either labeled data $\boldsymbol{x}_l$ or unlabeled data $\boldsymbol{x}_u$, and $D[p, q]$ is the KL divergence $D[p, q] = \sum_i p_i \log \frac{p_i}{q_i}$. Since no true label is required in the optimization above, VAT can be applied to both labeled and unlabeled data.

Fig. 1(a) illustrates the general training pipeline of AT and VAT. Due to the constrained optimization in Eqs. 1 and 2, the exact closed-form solution of the adversarial perturbation $\boldsymbol{r}_{\text{adv}}$ is intractable. Instead, fast approximation algorithms are

proposed to estimate $\boldsymbol{r}_{\text{adv}}$ iteratively. For AT, the adversarial perturbations can be approximated as:

$$\boldsymbol{r}_{\text{adv}} \approx \begin{cases} \epsilon \frac{\boldsymbol{g}}{\|\boldsymbol{g}\|_2} & L_2\text{-norm} \\ \epsilon \text{sign}(\boldsymbol{g}) & L_\infty\text{-norm}, \end{cases} \quad (3)$$

where $\boldsymbol{g} = \nabla_{\boldsymbol{x}_l} D\left[h(y|\boldsymbol{x}_l), p(y|\boldsymbol{x}_l, \boldsymbol{\theta})\right]$. And for VAT, the perturbation can approximated via the power iteration and estimated by:

$$\boldsymbol{r}_{\text{adv}} \approx \epsilon \frac{\boldsymbol{g}}{\|\boldsymbol{g}\|_2} \quad L_2\text{-norm}, \quad (4)$$

where $\boldsymbol{g} = \nabla_{\boldsymbol{r}} D\left[p(y|\boldsymbol{x}_*, \boldsymbol{\theta}), p(y|\boldsymbol{x}_* + \boldsymbol{r}, \boldsymbol{\theta})\right]$. For both algorithms, the backpropagation is needed to compute the additive perturbation $\boldsymbol{r}_{\text{adv}}$.

### B. Multiplicative Perturbations

In contrast to the additive perturbations exploited by AT and VAT, we introduce a new type of perturbations $\boldsymbol{z}$ that are multiplicative:

$$\boldsymbol{x}_{\text{xadv}} = \boldsymbol{x} \odot \boldsymbol{z}, \quad (5)$$

where $\boldsymbol{x} \in \mathbb{R}^P$ denotes an input image, $\boldsymbol{z} \in \{0, 1\}^P$ is a set of binary masks, and $\odot$ is the element-wise multiplication. We can interpret $\boldsymbol{x} \odot \boldsymbol{z}$ as an operation that attaches a binary random variable $z^j$ to each pixel $j$ of $\boldsymbol{x}$, for all $j \in \{1, 2, \cdots, P\}$. When $z^j = 0$, the corresponding pixel value is set to 0. Otherwise, the corresponding pixel value is retained without any changes. With the multiplicative perturbations, we have the following constrained optimization problem:

$$\mathcal{L}_{\text{xadv}}(\boldsymbol{x}, \boldsymbol{z}_{\text{xadv}}, \boldsymbol{\theta}) = D\left[p(y|\boldsymbol{x}, \boldsymbol{\theta}), p(y|\boldsymbol{x} \odot \boldsymbol{z}_{\text{xadv}}, \boldsymbol{\theta})\right] \quad (6a)$$
$$\text{with } \boldsymbol{z}_{\text{xadv}} = \underset{\boldsymbol{z}}{\arg\max} D\left[p(y|\boldsymbol{x}, \boldsymbol{\theta}), p(y|\boldsymbol{x} \odot \boldsymbol{z}, \boldsymbol{\theta})\right], \quad (6b)$$

where $D[p, q]$ adopts the cross entropy function for xAT, and the KL divergence for xVAT. To simplify the notation, we define

$$\Delta D(\boldsymbol{z}, \boldsymbol{x}, \boldsymbol{\theta}) = D\left[p(y|\boldsymbol{x}, \boldsymbol{\theta}), p(y|\boldsymbol{x} \odot \boldsymbol{z}, \boldsymbol{\theta})\right]. \quad (7)$$

In this formulation, $\boldsymbol{z}_{\text{xadv}} \in \{0, 1\}^P$ is optimized adversarially to maximize the divergence measure in Eq. 6b. This means that we wish $z^j$ to take value of 0 if zeroing out the corresponding pixel $j$ will make the prediction significantly different from the original prediction $p(y|\boldsymbol{x}, \boldsymbol{\theta})$. Apparently, a trivial solution of $\boldsymbol{z}$ is all 0s, which is likely to maximize the divergence measure in Eq. 6b, but is catastrophic to the optimization of model parameter $\boldsymbol{\theta}$ in Eq. 6a. To avoid this

detrimental solution, we augment Eq. 6b with the $L_0$-norm of $z$ to regularize the learning of $z$:

$$
\begin{aligned}
\boldsymbol{z}_{\text{xadv}} &= \arg\max_{\boldsymbol{z}} \Delta D(\boldsymbol{z}, \boldsymbol{x}, \boldsymbol{\theta}) + \lambda \|\boldsymbol{z}\|_0 \\
&= \arg\max_{\boldsymbol{z}} \Delta D(\boldsymbol{z}, \boldsymbol{x}, \boldsymbol{\theta}) + \lambda \sum_{j=1}^{P} 1_{[z^j \neq 0]}
\end{aligned}
\tag{8}
$$

where $1_{[c]}$ is an indicator function that outputs 1 if the condition $c$ is satisfied, and 0 otherwise. Consider two extreme cases. When $\boldsymbol{z} = \boldsymbol{0}$, the first term $\Delta D$ is likely to reach its maximum, but the second term $\|\boldsymbol{z}\|_0$ is minimized to 0, and thus $\boldsymbol{z} = \boldsymbol{0}$ is unlikely to maximize Eq. 8. On the other hand, when $\boldsymbol{z} = \boldsymbol{1}$, the first term $\Delta D$ is minimized to 0 and the second term $\|\boldsymbol{z}\|$ reaches its maximum, and thus $\boldsymbol{z} = \boldsymbol{1}$ is unlikely to maximize Eq. 8 either. Therefore, a good solution must lie in between these two extremes, where some elements of $\boldsymbol{z}$ are 0s and the remaining are 1s. As the training proceeds, the optimized $\boldsymbol{z}$ shall gradually become an adversarial mask that blocks salient regions of an image, leads to an unreliable prediction and therefore maximizes Eq. 8. Subsequently, this adversarial mask will regularize the trained DNN from Eq. 6a to be robust to this multiplicative perturbation. We will demonstrate this behavior when we present results.

*1) Stochastic Variational Optimization:* To optimize Eq. 8, we need to compute its gradient w.r.t. $\boldsymbol{z}$. Since $\boldsymbol{z} \in \{0,1\}^P$ is a set of binary random variables, both the first term and the second term of Eq. 8 are not differentiable. Hence, we resort to approximation algorithms to solve this binary optimization problem. We can use an inequality from stochastic variational optimization [7] to derive a lower-bound of Eq. 8. That is, given any function $\mathcal{F}(\boldsymbol{z})$ and any distribution $q(\boldsymbol{z})$, the following inequality holds:

$$
\max_{\boldsymbol{z}} \mathcal{F}(\boldsymbol{z}) \geq \mathbb{E}_{\boldsymbol{z} \sim q(\boldsymbol{z})} [\mathcal{F}(\boldsymbol{z})]
\tag{9}
$$

i.e., the maximum of a function is lower bounded by the expectation of the function.

Since $z^j, j \in \{1, 2, \cdots, P\}$ is a binary random variable, we assume $z^j$ is subject to a Bernoulli distribution with parameter $\pi^j \in [0,1]$, i.e. $q(z^j|\pi^j) = \text{Ber}\left(z^j; \pi^j\right)$. Thus, Eq. 8 can be lower bounded by its expectation:

$$
\boldsymbol{\pi}_{\text{xadv}} = \arg\max_{\boldsymbol{\pi}} \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{\pi})} \left[\Delta D(\boldsymbol{z}, \boldsymbol{x}, \boldsymbol{\theta})\right] + \lambda \sum_{j=1}^{P} \pi^j
\tag{10}
$$

Now the second term of Eq. 10 is differentiable w.r.t. the new model parameters $\boldsymbol{\pi}$. However, the first term is still problematic as the expectation over a large number of binary random variables $\boldsymbol{z}$ is intractable and so is its gradient. Therefore, further approximations are required.

*2) The Hard Concrete Gradient Estimator:* Thanksfully, this stochastic binary optimization problem has been investigated extensively in the literature. There exist a number of gradient estimators to this problem, such as REINFORCE [8], Gumble-Softmax [9, 10], REBAR [11], RELAX [12] and the hard concrete estimator [13]. We resort to the hard concrete

estimator to optimize Eq. 10 since it's straightforward to implement and demonstrates superior performance in our experiments. Specifically, the hard concrete gradient estimator employs a reparameterization trick to approximate the original optimization problem (10) by a close surrogate function:

$$
\begin{aligned}
\log \boldsymbol{\alpha}_{\text{xadv}} =& \arg\max_{\log \boldsymbol{\alpha}} \mathbb{E}_{\boldsymbol{u} \sim \mathcal{U}(0,1)} \left[\Delta D(g(f(\log \boldsymbol{\alpha}, \boldsymbol{u})), \boldsymbol{x}, \boldsymbol{\theta})\right] \\
&+ \lambda \sum_{j=1}^{P} \sigma\left(\log \alpha^j - \beta \log \frac{-\gamma}{\zeta}\right)
\end{aligned}
\tag{11}
$$

with

$$
f(\log \boldsymbol{\alpha}, \boldsymbol{u}) = \sigma\left((\log \boldsymbol{u} - \log(1-\boldsymbol{u}) + \log \boldsymbol{\alpha})/\beta\right)(\zeta - \gamma) + \gamma,
$$
$$
g(\cdot) = \min(1, \max(0, \cdot)),
$$

where $\mathcal{U}(0,1)$ denotes the uniform distribution in the range of $[0,1]$, $\sigma(t) = 1/(1 + \exp(-t))$ is the sigmoid function, and $\beta = 2/3$, $\gamma = -0.1$, and $\zeta = 1.1$ are the typical values of the hard concrete distribution. With this reparameterization, the surrogate function (11) is now differentiable w.r.t. its parameters $\log \boldsymbol{\alpha} \in \mathbb{R}^P$. For more details on the hard concrete gradient estimator, we refer the readers to [13].

Once we have an optimized $\log \boldsymbol{\alpha}_{\text{xadv}}$, we can sample an adversarial mask from the hard concrete distribution $q(\boldsymbol{z}|\log \boldsymbol{\alpha}_{\text{xadv}})$ by

$$
\hat{\boldsymbol{z}}_{\text{xadv}} = g(f(\log \boldsymbol{\alpha}_{\text{xadv}}, \boldsymbol{u})), \qquad \boldsymbol{u} \sim \mathcal{U}(0,1).
\tag{12}
$$

We can then optimize model parameter $\boldsymbol{\theta}$ by minimizing the following regularized empirical risk over all training examples:

$$
\begin{aligned}
\mathcal{L}(\boldsymbol{\theta}) =& \frac{1}{N_l} \sum_{(\boldsymbol{x}_l, y_l) \in \mathcal{D}_l} \mathcal{L}_{\text{ce}}(h(\boldsymbol{x}_l, \boldsymbol{\theta}), y_l) \\
&+ \eta \frac{1}{N_l + N_u} \sum_{\boldsymbol{x} \in \{\mathcal{D}_l, \mathcal{D}_u\}} \mathcal{L}_{\text{xadv}}(\boldsymbol{x}, \hat{\boldsymbol{z}}_{\text{xadv}}, \boldsymbol{\theta}),
\end{aligned}
\tag{13}
$$

where $h(\boldsymbol{x}_l, \boldsymbol{\theta})$ denotes the prediction of a classifier, parameterized by $\boldsymbol{\theta}$, $\mathcal{L}_{\text{ce}}$ is the cross entropy loss over labeled training examples, and $\eta$ is a regularization hyperparameter that balances the cross entropy loss $\mathcal{L}_{\text{ce}}$ of labeled data and the adversarial training loss $\mathcal{L}_{\text{xadv}}$ of all training examples.

*3) Transductive vs. Inductive Training:* The optimization of the final loss function (13) can be implemented transductively or inductively. For a transductive implementation, we set $\log \boldsymbol{\alpha} \in \mathbb{R}^P$ as model parameters to be optimized for each input example $\boldsymbol{x} \in \{\mathcal{D}_l, \mathcal{D}_u\}$. The learned $\log \boldsymbol{\alpha}$ is then used to generate a sparse adversarial mask $\hat{\boldsymbol{z}}_{\text{xadv}}$, and subsequently the training proceeds to evaluate the final loss (13). However, this approach can't generate masks for images never seen during training. A more appealing approach is the inductive implementation, which employs a generator to produce a sparse adversarial mask for any input example $\boldsymbol{x}$. To formulate the inductive training, we model the generator as a neural network, parameterized by $\boldsymbol{\theta}_g$, and given an input $\boldsymbol{x}$ we define its output as $\log \boldsymbol{\alpha} = G(\boldsymbol{x}, \boldsymbol{\theta}_g)$, which is used subsequently to sample a sparse adversarial mask $\hat{\boldsymbol{z}}_{\text{xadv}}$. The pipeline of the two implementations is illustrated in Fig. 2, which can be optimized end-to-end.
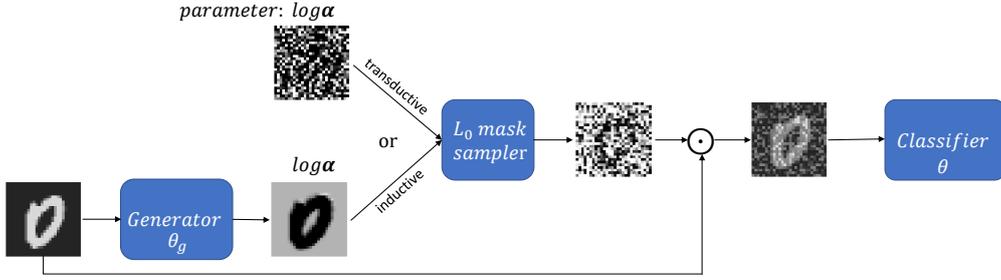
Fig. 2. The pipeline of the transductive and inductive implementations of multiplicative adversarial training.

One advantage of xAT/xVAT over AT/VAT is that due to the hard concrete reparameterization, we can update (i) classifier parameter $\theta$, and (ii) the hard concrete parameter $\log \alpha$ (for transductive training) or the generator parameter $\theta_g$ (for inductive training) simultaneously in one step (See Fig. 2). This is not true for AT and VAT since both of them rely on backpropagation to optimize additive perturbations, while backpropagating through the perturbations produced by backpropagation is computationally unstable. Therefore, both AT and VAT resort to optimizing additive perturbations $r_{\text{adv}}$ and classifier parameter $\theta$ alternatively in two steps. Because of this advantage, xAT and xVAT are computationally more efficient than their additive counterparts. We will demonstrate this when we present results (Table III).
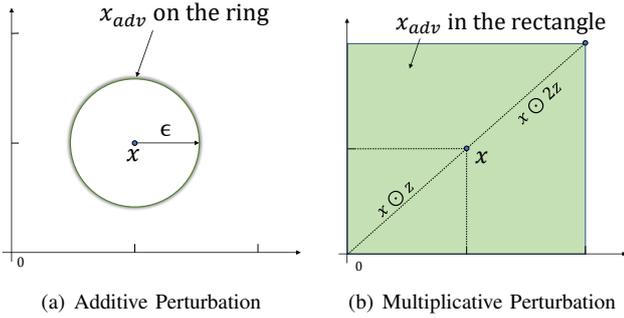


(a) Additive Perturbation  (b) Multiplicative Perturbation

Fig. 3. The effect of $\epsilon$ on different perturbations. (a) shows that the additive perturbations are on the surface of a ball with the radius $\epsilon$. (b) demonstrates that our multiplicative perturbations are distributed within the rectangle.

*4) Shrink or Expand Multiplicative Perturbations:* It is worth mentioning one potential issue of multiplicative perturbations. Consider the additive adversarial example $x + r_{\text{adv}}$ with $r_{\text{adv}} \approx \epsilon \frac{g}{\|g\|_p}$ (where $p = 2$ or $\infty$), and the multiplicative adversarial example $x \odot z$ with $z \in \{0, 1\}^P$. Apparently, these two types of perturbations have very distinct geometric interpretations w.r.t. input example $x$. For additive perturbations, the approximation $r_{\text{adv}} \approx \epsilon \frac{g}{\|g\|_p}$ indicates that all additive perturbations are generated exactly on the surface of a ball (or a box) centered at $x$ with a radius of $\epsilon$, and no additive perturbations would appear inside the ball (see Fig. 3(a)). On the other hand, the multiplicative perturbations do not have this characteristic. Since $z^j$ is between 0 and 1, the admissible multiplicative perturbations can appear inside a box rather than only on the surface of a box (see Fig. 3(b)). Note that the following inequality always holds:

$$\|x \odot z\| \le \|x\|, \text{ with } z \in [0, 1]^P. \qquad (14)$$

What does this mean is that the masked image $x \odot z$ is always darker than the original image $x$, which may be undesirable in some cases where the contrast of an image is low. To tackle this potential issue, we introduce a new parameter $\epsilon$ in the multiplicative formulation: $x \odot \epsilon z$. We only consider $\epsilon = 1$ or 2 in our experiments, while other positive values of $\epsilon$ are admissible. When $\epsilon = 1$, the masked images can only be darker than the original images, while when $\epsilon = 2$ the masked images can be darker or brighter as illustrated in Fig. 3(b). Empirically, we find that $\epsilon = 1$ works reasonable well on all benchmark datasets in our experiments. We therefore don't tune it any further.

Comparing Fig. 3(a) and Fig. 3(b), it is clear that the solution space of multiplicative perturbations is much larger than that of additive perturbations. As we will show in the experiments, due to this difference, the trained DNNs from both perturbations demonstrate very distinct weight distributions.

## III. RELATED WORK

To improve the generalization of trained DNNs on unseen data examples, a variety of regularization techniques have been proposed in the past few years. Traditional regularization techniques impose an extra term to prevent the model from overfitting to a small set training examples. Recenlty, a family of regularization techniques called consistency regularization has been proposed that encourages classifiers to yield consistent predictions on data samples and their perturbed versions. For example, the $\Pi$-Model [14] and STP [15] incorporate the following consistency loss

$$\|p(y|\text{Augment}_1(x), \theta) - p(y|\text{Augment}_2(x), \theta)\|_2^2 \qquad (15)$$

as a regularization for robust training of DNNs, where Augment$(x)$ represents a stochastically perturbed version of $x$ so that the two predictions in Eq. 15 are noisy, and minimizing their difference may improve the robustness.

Similarly, Temporal Ensembling [14], Mean Teacher [16] and fast-SWA [17] generate stochastic outputs by tracking the exponential moving average (EMA) of the past predictions and weights. AT and VAT generate small but worse-case perturbations as data augmentation. All of these approaches encourage classifiers to produce consistent predictions on different perturbed examples. More recently, by combining the dropout and adversarial training, Park et al. [18] propose
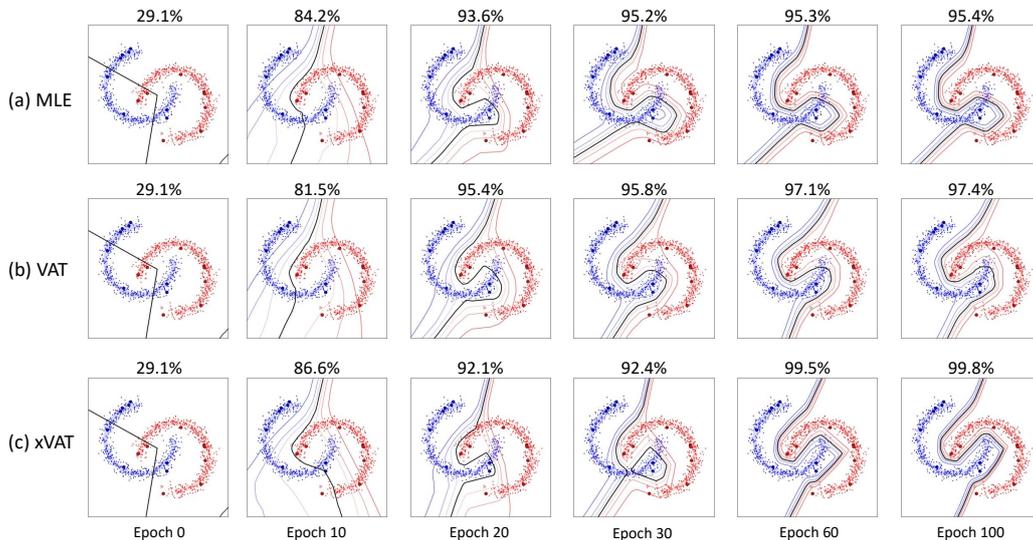
Fig. 4. The evolution of the decision boundaries of MLE, VAT and xVAT on the synthetic "moons" dataset.

the adversarial dropout, which embeds an adversarial dropout layer to DNNs and optimizes the dropout mask adversarially to maximize the loss of DNNs. This work has a similar high-level idea to ours. However, their method works on much smaller last FC layer, and employs an expensive integer programming to solve the optimization problem. Because of these limitations, their technique cannot be applied to earlier but larger layers, while our xAT and xVAT operate at input layer with an $L_0$-regularized binary optimization, which can be optimized efficiently via the hard concrete gradient estimator. Our experiments verify that xVAT not only demonstrates the state-of-the-art accuracies but also is about 30% faster than VAT, and therefore is much more scalable.

## IV. EXPERIMENTAL RESULTS

We validate our xAT/xVAT on multiple datasets with different network architectures for supervised learning and semi-supervised learning. Specifically, we illustrate how xVAT learns on a synthetic "moons" dataset [4]. We also demonstrate xAT and xVAT on four established image classification benchmarks: MNIST [19], SVHN [20], CIFAR-10 and CIFAR-100 [21]. Our main baselines are AT [3] and VAT [6] since our work is primarily to investigate the effectiveness of multiplicative perturbations *vs.* additive perturbations. We also compare xAT/xVAT with many other algorithms in terms of classification accuracy. For a fair comparison, our implementation and experimental setup [1] closely follows that of VAT [2].

### A. Synthetic Dataset

We first demonstrate the performance of xVAT on a synthetic "Moons" dataset [4]. The dataset contains 16 labeled training data points and 1000 test data points, uniformly distributed in two moon-shaped clusters for binary classification. We randomly select 500 test points as unlabeled training examples and use all 1000 test points for evaluation.

We compare the results of Maximum Likelihood Estimation (MLE) using the cross-entropy loss, VAT and xVAT. MLE only uses 16 labeled data points for supervised training, while VAT and xVAT use 16 labeled and 500 unlabeled data points for semi-supervised training. Fig. 4 illustrates the evolution of the decision boundaries of MLE, VAT, xVAT on the "moons" dataset. As can be seen, xVAT yields a similar or slightly better decision boundary than that of VAT. Furthermore, both VAT and xVAT lean the decision boundaries that respect the underlying data manifold, demonstrating the effectiveness of multiplicative perturbations on this synthetic dataset.

### B. Image Classification Benchmarks

We next evaluate xAT and xVAT on four established image classification benchmarks. Detailed description of the benchmarks, network architectures and experimental setup for reproducibility can be found in supplementary material and our open source implementation.

*1) Semi-supervised Learning:* For MNIST, we use an MLP with four hidden layers of 1200, 600, 300 and 150 units, respectively, which is the same architecture used in VAT [4, 6]. For SVHN, CIFAR-10 and CIFAR-100, the network architecture is a 13-layer CNN, the same as the one used in [6, 14, 18]. It is worth mentioning that for xAT and xVAT we normalize the MNIST images to $[-0.5, 0.5]$ rather than $[0, 1]$. This is because multiplicative perturbations have no effect on zero-valued pixels, while the MNIST images happen to have zero-valued black background, and therefore the $[-0.5, 0.5]$ normalization allows xAT/xVAT to generate valid multiplicative perturbations on the background of MNIST. On the other datasets, we adopt the same normalization approaches as in VAT. For the transductive implementation, we initialize the parameter $\log \alpha$ for each image with samples from a unit Gaussian distribution $\mathcal{N}(0, 1)$. For the inductive implementation, we use an one-layer CNN with one $3 \times 3$ filter as the generator. Interestingly, such a simple generator is sufficient to yield competitive results for all our experiments.

[1] https://github.com/sndnyang/xvat/
[2] https://github.com/takerum/vat_tf/

TABLE I
TEST ACCURACIES OF SEMI-SUPERVISED LEARNING ON MNIST, SVHN AND CIFAR-10. THE RESULTS ARE AVERAGED OVER 5 RUNS.

| Method | Test Accuracy (%) | | |
| | MNIST $N_l$=100 | SVHN $N_l$=1000 | CIFAR-10 $N_l$=4000 |
| --- | --- | --- | --- |
| GAN with feature match [22] | **99.07** | 91.89 | 81.37 |
| CatGAN [23] | 98.09 | - | 80.42 |
| Ladder Networks [24] | 98.94 | - | 79.60 |
| Π-model [14] | - | 94.57 | 83.45 |
| Mean Teacher [16] | - | **94.79** | 82.26 |
| VAT [6] | 98.64 | 94.23 | 85.18 |
| xVAT (Transductive) | 98.02 | 93.99 | 85.82 |
| xVAT (Inductive) | 97.82 | 94.22 | **86.59** |

TABLE II
TEST ACCURACIES OF SEMI-SUPERVISED LEARNING ALGORITHMS ON CIFAR-100. THE RESULTS ARE AVERAGED OVER 5 RUNS.

| Method | Test Accuracy (%) $N_l$=10000 |
| --- | --- |
| Supervised [14] | 55.44 |
| Π-model [14] | 60.81 |
| Temporal ensembling [14] | 61.35 |
| VAT [6] | 59.71 |
| xVAT (Transductive) | 61.30 |
| xVAT (Inductive) | **61.76** |

TABLE III
THE TRAINING SPEEDS OF VAT AND xVAT ON THE FOUR BENCHMARK DATASETS. THE RESULTS ARE AVERAGED OVER 5 RUNS.

| Method | Seconds per epoch | | | |
| | MNIST | SVHN | CIFAR-10 | CIFAR-100 |
| --- | --- | --- | --- | --- |
| VAT (ours)* | 4.31 | 54.3 | 51.3 | 51.5 |
| xVAT (Transductive) | 4.54 | 36.6 | 34.1 | 39.3 |
| xVAT (Inductive) | 4.33 | 35.7 | 33.6 | 34.4 |

*For a fair comparison, we implement VAT in PyTorch and achieve similar accuracies as the official VAT implementation.

Tables I and II summarize the classification accuracies of xVAT on the four benchmark datasets as compared to the state-of-the-art algorithms. As we can see, xVAT achieves very competitive and sometimes even better accuracies than the state-of-the-arts, such as the Π-model [14] and VAT. Furthermore, the inductive xVAT outperforms the transductive xVAT on 3 out of 4 benchmarks. This is likely because the per-image parameter $\log \alpha$ is more difficult to optimize than that of the shared generator parameter $\theta_g$. Table III compares the per-epoch training speeds of VAT and xVAT. xVAT (both transductive and Inductive) is about 30% faster than VAT. This is because the generator of xVAT is significantly shallower than the CNN classifier, and therefore the multiplicative perturbations can be generated much more efficiently than the backpropagation-based additive perturbations. In general, VAT requires at least two pairs of forward and backward propagations per iteration, while xVAT only needs one pair of forward and backward propagation, plus some extra time to update the parameter of generator. In our experiments, since an one-layer CNN generator is sufficient, which is significantly cheaper than the 13-layer CNN classifier, the cost of updating the generator is largely negligible, and therefore our xVAT is computationally more efficient than VAT. Because inductive xAT/xVAT in general outperform their transductive ones, in the experiments that follow we use inductive xAT/xVAT unless noted otherwise.

TABLE IV
TEST ACCURACIES OF SUPERVISED LEARNING ON MNIST. THE RESULTS ARE AVERAGED OVER 5 RUNS.

| Method | Test Accuracy (%) |
| --- | --- |
| Dropout [1] | 98.95 |
| Concrete Dropout [25] | 98.60 |
| Ladder networks [24] | **99.43** |
| Baseline (MLE) [6] | 98.89 |
| AT, $L_\infty$ [6] | 99.21 |
| AT, $L_2$ [6] | 99.29 |
| VAT [6] | 99.36 |
| xAT (Inductive) | 99.18 |
| xVAT (Inductive) | 99.08 |

TABLE V
TEST ACCURACIES OF SUPERVISED LEARNING ON CIFAR-10 AND CIFAR-100. THE RESULTS ARE AVERAGED OVER 5 RUNS.

| Method | Test Accuracy (%) | |
| | CIFAR-10 | CIFAR-100 |
| --- | --- | --- |
| Baseline (MLE) [14] | 93.24 | 73.58 |
| Π-model [14] | **94.44** | 73.68 |
| Temporal ensembling [14] | 94.40 | 73.70 |
| AT, $L_\infty$ (ours)* | 93.90 | 74.04 |
| VAT [6] | 94.19 | 75.02 |
| xAT (Inductive) | 93.70 | 74.62 |
| xVAT (Inductive) | 93.88 | **75.30** |

*No reported results. For a fair comparison, we implement AT in PyTorch and verify its accuracies on MNIST, and report our results in the table.

*2) Supervised Learning:* We next evaluate the performance of xAT and xVAT in supervised learning on MNIST, CIFAR-10 and CIFAR-100. The same network architectures are used as in the experiments of semi-supervised learning. We compare the results of xAT and xVAT with the state-of-the-art algorithms in Tables IV and V. It is demonstrated that our xAT and xVAT achieve very competitive accuracies to the other competing methods, demonstrating the effectiveness of multiplicative perturbations in supervised learning.

*C. Visualization of Multiplicative Perturbations*

To understand the effectiveness of xVAT in generating multiplicative perturbations, we visualize the evolution of $\log \alpha$ and $x_{\text{xadv}}$ on example images from MNIST, SVHN and CIFAR-10, with the results shown in Fig. 5(a). As we can see, at the beginning of the xVAT training, the output of sparse mask generator $\log \alpha$ isn't very informative as it generates multiplicative perturbations almost uniformly over the entire images. As the training proceeds and the sparse mask generator is trained adversarially, the mask learns to block the salient regions of the input images and leads to more effective adversarial examples. This can be observed from the learned $\log \alpha$, which gets more interpretable over training epochs.
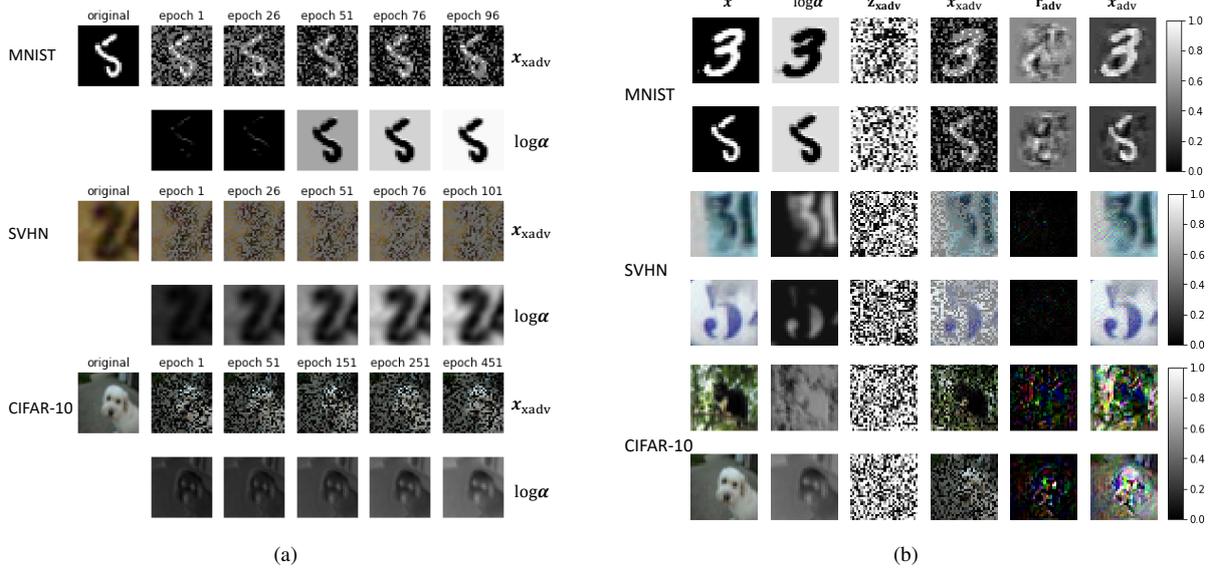
Fig. 5. Visualization of multiplicative perturbations and additive perturbations from xVAT and VAT. (a) The evolution of $\log \boldsymbol{\alpha}$ and $\boldsymbol{x}_{\text{xadv}}$ during the training of xVAT on benchmark datasets. (b) Comparison of multiplicative and additive perturbations on example images from benchmark datasets.
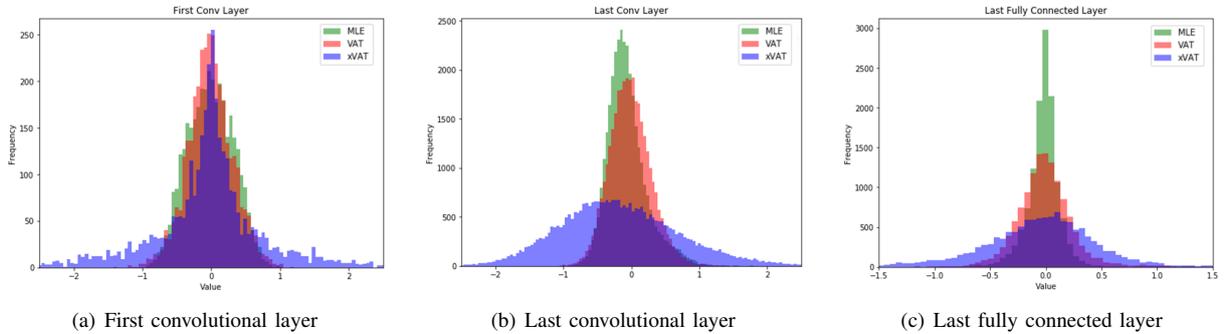


Fig. 6. Histograms of the classifier weights learned by MLE, VAT and xVAT on CIFAR-100. The histograms are computed from different CNN layers.

To compare the multiplicative perturbations with the additive perturbations, Fig. 5(b) illustrates their effects on example images from MNIST, SVHN and CIFAR-10. As we can see, visually these two types of perturbations are very different: the multiplicative perturbations are more perceptible than the additive counterparts, but the former are more interpretable than the latter as learned $\log \boldsymbol{\alpha}$s clearly discover the salient regions of input images. It is worth mentioning that the more perceptible multiplicative perturbation is not a disadvantage of xAT/xVAT since our work is not to propose new adversarial attack algorithms but to investigate new adversarial training techniques that can improve the robustness of DNNs. Evidently, these results demonstrate that multiplicative perturbations are *at least* as effective as additive perturbations at identifying the blind-spots of DNNs and training on these perturbations can lead more robust DNNs.

### D. Visualization of Weight Distributions

To investigate the differences of trained classifiers from different algorithms, we compare the weights of models learned by MLE, VAT and xVAT on CIFAR-100. Fig. 6 shows the histograms of model weights from three different

CNN layers. It is demonstrated that xVAT learns a classifier whose weights have a much higher variance than those learned by VAT and MLE. In other words, xVAT learns a denser classifier from multiplicative perturbations with more non-zero weights than VAT and MLE. As we have demonstrated in Fig. 3, the multiplicative perturbations reside inside a rectangle around an input $\boldsymbol{x}$, while the additive perturbations can only lie on the ring surrounding $\boldsymbol{x}$. As a result, there are more varieties among multiplicative perturbations than their additive counterparts, and therefore the adversarially trained DNNs need more capacities (active neurons) to against multiplicative perturbations.

## V. CONCLUSION

In this paper we propose a new type of perturbations that is multiplicative. Compared to the additive perturbations exploited by AT and VAT, the multiplicative perturbations are more perceptible and interpretable. We show that these multiplicative perturbations can be optimized via an $L_0$-norm regularized objective and implemented transductively and inductively. Thanks to the hard concrete reparameterization, the resulting algorithms xAT and xVAT are computationally more

efficient than their additive counterparts. Extensive experiments on synthetic and four established image classification benchmarks demonstrate that xAT and xVAT match or outperform the state-of-the-art algorithms while being about 30% faster. Visualization of xAT and xVAT demonstrates a stark contract to their additive counterparts.

As for future extensions, we are interested in investigating the interplay between xVAT and $L_0$-norm based network sparsification [26]. Both algorithms are developed via the $L_0$ regularizations, while xVAT learns a dense network and $L_0$-ARM [26] prunes redundant neurons. The combination of them might lead a network that is both robust and efficient.

## VI. ACKNOWLEDGMENT

## REFERENCES

[1] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, 2014.

[2] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *International Conference on Learning Representations (ICLR)*, 2014.

[3] I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *International Conference on Learning Representations*, 2015.

[4] T. Miyato, S. Maeda, M. Koyama, K. Nakae, and S. Ishii, "Distributional smoothing by virtual adversarial examples," in *International Conference on Learning Representations (ICLR)*, 2016.

[5] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "Mixup: Beyond empirical risk minimization," in *International Conference on Learning Representations (ICLR)*, 2018.

[6] T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii, "Virtual adversarial training: a regularization method for supervised and semi-supervised learning," *IEEE transactions on PAMI*, 2018.

[7] T. Bird, J. Kunze, and D. Barber, "Stochastic variational optimization," *arXiv preprint arXiv:1809.04855*, 2018.

[8] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine Learning*, 1992.

[9] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," in *International Conference on Learning Representations (ICLR)*, 2017.

[10] C. J. Maddison, A. Mnih, and Y. W. Teh, "The concrete distribution: A continuous relaxation of discrete random variables," in *International Conference on Learning Representations (ICLR)*, 2017.

[11] G. Tucker, A. Mnih, C. J. Maddison, J. Lawson, and J. Sohl-Dickstein, "Rebar: Low-variance, unbiased gradient estimates for discrete latent variable models," in *Advances in Neural Information Processing Systems (NIPS)*, 2017.

[12] W. Grathwohl, D. Choi, Y. Wu, G. Roeder, and D. Duvenaud, "Backpropagation through the void: Optimizing control variates for black-box gradient estimation," in *International Conference on Learning Representations (ICLR)*, 2018.

[13] C. Louizos, M. Welling, and D. P. Kingma, "Learning sparse neural networks through $l_0$ regularization," in *International Conference on Learning Representations (ICLR)*, 2018.

[14] L. Samuli and A. Timo, "Temporal ensembling for semi-supervised learning," in *International Conference on Learning Representations (ICLR)*, 2017.

[15] M. Sajjadi, M. Javanmardi, and T. Tasdizen, "Regularization with stochastic transformations and perturbations for deep semi-supervised learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.

[16] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[17] B. Athiwaratkun, M. Finzi, P. Izmailov, and A. G. Wilson, "There are many consistent explanations of unlabeled data: Why you should average," in *International Conference on Learning Representations*, 2019.

[18] S. Park, J. Park, S.-J. Shin, and I.-C. Moon, "Adversarial dropout for supervised and semi-supervised learning," in *AAAI Conference on Artificial Intelligence*, 2018.

[19] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," in *Proceedings of the IEEE*, 1998.

[20] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and N. andrew Y, "Reading digits in natural images with unsupervised feature learning," 2011.

[21] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Tech. Rep., 2009.

[22] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *Advances in neural information processing systems (NeurIPS)*, 2016.

[23] J. T. Springenberg, "Unsupervised and semi-supervised learning with categorical generative adversarial networks," in *International Conference on Learning Representations (ICLR)*, 2015.

[24] A. Rasmus, M. Berglund, M. Honkala, H. Valpola, and T. Raiko, "Semi-supervised learning with ladder networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.

[25] Y. Gal, J. Hron, and A. Kendall, "Concrete dropout," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[26] Y. Li and S. Ji, "L0-ARM: Network sparsification via stochastic binary optimization," in *The European Conference on Machine Learning (ECML)*, 2019.

[27] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2015.

## VII. APPENDIX

### A. Image Classification Benchmarks

The image classification benchmarks used in our experiments are described below:

1) MNIST [19] is a gray-scale image dataset containing 60,000 training images and 10,000 test images of the size $28 \times 28$ for 10 handwritten digits classification.
2) SVHN [20] is a street view house number dataset containing 73,257 training images and 26,032 test images classified into 10 classes representing digits. Each image may contain multiple real-world house number digits, and the task is to classify the center-most digit. These are RGB Images of the size $32 \times 32$.
3) CIFAR-10 [21] contains 10 classes of RGB images of the size $32 \times 32$, in which 50,000 images are for training and 10,000 images are for test.
4) CIFAR-100 [21] also has 60,000 RGB images of the size $32 \times 32$, except that it contains 100 classes with 500 training images and 100 test images per class.

### B. Experimental Setup

Our xAT and xVAT contains two sets of model parameters: (i) classifier parameter $\boldsymbol{\theta}$, and (ii) per-image $\log \boldsymbol{\alpha}$ parameter as in the transductive implementation or generator parameter $\boldsymbol{\theta}_g$ as in the inductive implementation. For the transductive implementation, we optimize $\log \boldsymbol{\alpha}$ by SGD with a learning rate of 0.001. For the inductive implementation, we found that one-layer CNN with one $3 \times 3$ filter works well for all our experiments, and we optimize the generator by using Adam [27] with a fixed learning rate of $1e-6$ or $2e-6$.

As for classifier parameter $\boldsymbol{\theta}$, different network architectures are used for different image classification benchmarks. In all our experiments, we use $\lambda = 1, \eta = 1$ for xVAT, and $\lambda = 1, \eta = 0.5$ for xAT.

*a) MNIST:* We use an MLP with four hidden layers of 1200, 600, 300 and 150 units, respectively, for the MNIST experiments. The same MLP architecture is used for supervised training and semi-supervised training as in VAT [4, 6]. The input images are linearly normalized to the range of $[-0.5, 0.5]$. We use the Adam optimizer [27] with an initial learning rate of 0.002, which is decayed by 0.9 at every 500 iterations. We train the classifier for 50,000 iterations, and at each iteration a mini-batch of 100 labeled images and 250 unlabeled images is used for training.

*b) CNN Architecture for SVHN, CIFAR-10/100:* On these three benchmarks, for the purpose of fair comparisons, we use the CNN architecture shown in Table VI, which is the same as the one used in [6, 14, 18].

We follow the same normalization schemes of VAT on these three benchmarks. We normalize the SVHN images to the range of $[0, 1]$. For CIFAR-10 and CIFAR-100, we normalize

TABLE VI
THE CNN ARCHITECTURE USED ON SVHN, CIFAR-10, AND CIFAR-100.

| Name | Description |
|---|---|
| Input | Image |
| Conv1a | 128 filters, $3 \times 3$, Pad='same', LReLU($\alpha = 0.1$) |
| Conv1b | 128 filters, $3 \times 3$, Pad='same', LReLU($\alpha = 0.1$) |
| Conv1c | 128 filters, $3 \times 3$, Pad='same', LReLU($\alpha = 0.1$) |
| Pool1 | Maxpool $2 \times 2$ pixels, Stride 2 |
| Drop1 | Dropout, p = 0.5 |
| Conv2a | 256 filters, $3 \times 3$, Pad='same', LReLU($\alpha = 0.1$) |
| Conv2b | 256 filters, $3 \times 3$, Pad='same', LReLU($\alpha = 0.1$) |
| Conv2c | 256 filters, $3 \times 3$, Pad='same', LReLU($\alpha = 0.1$) |
| Pool2 | Maxpool $2 \times 2$ pixels, Stride 2 |
| Drop2 | Dropout, p = 0.5 |
| Conv3a | 512 filters, $3 \times 3$, Pad='valid', LReLU($\alpha = 0.1$) |
| Conv3b | 256 filters, $1 \times 1$, Pad='same', LReLU($\alpha = 0.1$) |
| Conv3c | 128 filters, $1 \times 1$, Pad='same', LReLU($\alpha = 0.1$) |
| Pool3 | Global average pool, $6 \times 6 \to 1 \times 1$ |
| Dense | Fully connected $128 \to 10$ |
| Output | Softmax |

the images with ZCA, based on training set statistics, the same as in VAT [4, 6].

In semi-supervised learning, we train the classifier for 50000 iterations on SVHN, 200000 iterations on CIFAR-10, and 120000 iterations on CIFAR-100. At each iteration, we sample a mini-batch of 32 labeled images and 128 unlabeled images. We use the Adam optimizer with an initial learning rate of 0.001 and decay the learning rate linearly at the $[35000, 180000, 80000]$-th iterations.

For supervised learning on CIFAR-10 and CIFAR-100, we train the CNN classifier for 300 epochs with a batch size of 100. We again use Adam with an initial learning rate of 0.003 for xAT and 0.001 for xVAT. We linearly decay the learning rate at the 2nd half of training.

### C. Hyperparameter Tuning

Fig. 7 shows the results of hyperparameter tuning of $\lambda$ and $\eta$ on the MNIST dataset. As we can see, $\lambda = 1$ and $\eta = 1$ yields very competitive classification accuracies. Similar results are observed on the other benchmarks. We therefore use $\lambda = 1$ and $\eta = 1$ in all our semi-supervised learning tasks.
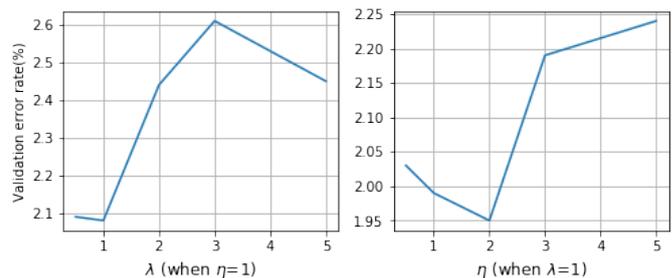


Fig. 7. Hyperparameters tuning of $\lambda$ and $\eta$ on the MNIST dataset.