# Non-Determinism in Neural Networks for Adversarial Robustness

#### Daanish Ali Khan

Department of Computer Science Carnegie Mellon University Pittsburgh, PA 15213 malikhan@andrew.cmu.edu

## Ninghao Sha

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213 nsha@andrew.cmu.edu

## Abelino Jimenez

Electrical and Computer Engineering Carnegie Mellon University Pittsburgh, PA 15213 abjimenez@cmu.edu

## Linhong Li

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213 linhongl@andrew.cmu.edu

#### Zhuoran (Oliver) Liu

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213 zhuoranl@andrew.cmu.edu

## Bhiksha Raj

Language Technologies Institute Carnegie Mellon University Pittsburgh, PA 15213 bhiksha@cs.cmu.edu

## Rita Singh

Language Technologies Institute Carnegie Mellon University Pittsburgh, PA 15213 rsingh@cs.cmu.edu

## **Abstract**

Recent breakthroughs in the field of deep learning have led to advancements in a broad spectrum of tasks in computer vision, audio processing, natural language processing and other areas. In most instances where these tasks are deployed in real-world scenarios, the models used in them have been shown to be susceptible to adversarial attacks, making it imperative for us to address the challenge of their adversarial robustness. Existing techniques for adversarial robustness fall into three broad categories: defensive distillation techniques, adversarial training techniques, and randomized or non-deterministic model based techniques. In this paper, we propose a novel neural network paradigm that falls under the category of randomized models for adversarial robustness, but differs from all existing techniques under this category in that it models each parameter of the network as a statistical distribution with learnable parameters. We show experimentally that this framework is highly robust to a variety of white-box and black-box adversarial attacks, while preserving the task-specific performance of the traditional neural network model.

## 1 Introduction

Neural systems are currently used in a broad spectrum of complex classification tasks, such as object recognition, speech processing, text generation etc. Many are deployed in large-scale tasks that are critical to human well-being and safety, such as biometric access points, medical assessments and self-driving cars. The systems themselves, however, are currently largely unprotected against malicious adversarial attacks – the presentation of inputs that have been purposely crafted to make the systems behave in incorrect ways [1] (Figure.1). Motivated largely by the desire to expose these vulnerabilities, a significant body of scientific literature has arisen in recent times on increasingly sophisticated techniques to generate adversarial instances – inputs that may fool machine learning systems [1, 2, 3].

With the increasing ubiquity of deep learning systems in the real world, the task of designing network architectures and learning paradigms that are robust to adversarial attacks is now recognized to be of paramount importance, and not surprisingly many solution approaches have been proposed in the literature [4, 5, 6]. Adversarial training attempts to adjust classifier decision boundaries away from adversarial instances by including the latter in the training data [7]. Distillation-based methods "distil" the trained networks into secondary networks to minimize their sensitivity to adversarial modifications of the input [8, 9]. Projection [10] and reconstruction methods [4] attempt to project down (and possibly reconstruct) inputs prior to feeding them to the system, in order to eliminate adversarial modifications. Randomization-based methods add noise or other random transformations [11] to the input to mask out adversarial modifications [12]. All of these methods assume the classification network itself to be deterministic.

In this paper, we propose a novel and alternate route to adversarial robustness. In our approach the *parameters of the network are themselves stochastic*, having a statistical distribution with learnable parameters. Inference on the network too is stochastic. The premise behind our model is that the randomness in the model confounds the ability of the adversary to determine the minimal change of the input required to fool it. Randomness during inference also increases the probability of avoiding the increased-variance adversarial inputs that result.

We use a modified version of the Stochastic Delta Rule (SDR)[13] to implement our stochastic models, employing a novel reparameterization trick to learn the distributions for the network parameters. We show experimentally that this framework does indeed provide protection against a variety of adversarial attacks in which other defences fail, while preserving the task-specific performance of the traditional neural network model.



Figure 1: Examples of adversarial attaks on a variety of application domains. From left to right: Adversarial physical-world attack on stop signs [14], face-recognition systems [15], and medical diagnosis systems [16]. The adversarilly designed patches on the stop sign, spectacles on the face and noise in the x-ray all cause them to be misrecognized.

## 2 Attacks: The Adversarial Threat Model

We briefly discuss existing approaches to generate adversarial samples, to set the background for the discussion of current state-of-the-art defense techniques and our own proposal.

The goal of the adversary is generally to generate examples that are perceptually indistinguishable from authentic ("clean") inputs, but are incorrectly classified by the model. The adversary may either choose to generate inputs that produce a specific (bogus) output from the classifier (*targeted attack*), or modify a clean input such that it is classified as not belonging to its true class, without explicitly considering what it may be classified as instead (*untargeted attack*). We primarily consider the latter in this paper, although the proposed approach should generalize to targeted attacks as well.

Adversarial examples can be generated in either a white-box or a black-box setting. In the white-box scenario, the adversary has full access to the target model's architecture and gradients. Adversarial inputs are generally obtained by minimally perturbing clean inputs such that they now produce bogus outputs [1]. The perturbations are computed using variations of gradient descent. After a forward pass, gradients of an adversary-defined objective are back-propagated onto the clean input, revealing the adversarial perturbations that would confuse the model once added to it. Adversarial perturbations may be generated either by taking a single step along the gradient (*one-step* methods) [1] or taking steps iteratively until some stopping criterion is met (*iterative* methods) [7]. On the other hand, black-box adversarial attacks craft adversarial examples without any internal knowledge of the target network, which makes them much more applicable in the real-world setting. Here the general approach is to probe the classifier with inputs to obtain input-output pairs. These are used to learn how to generate adversarial samples [17].

In our experiments, we choose to evaluate model robustness by classification accuracy under two white-box attacks, the Fast Gradient Sign Method (FGSM) and DeepFool Attack, and one black-box attack named LocalSearch, which we briefly review below.

## 2.1 White-box one-step attack example: Fast Gradient Sign Method

Define the loss  $L(X,y^{target})$  as a function of input X and its target label  $y^{target}$ , which regular, non-adversarial training tries to minimize. To produce adversarial instances the adversary instead increases  $L(X,y^{target})$ , tweaking the input such that the model is less likely to classify it correctly. The **fast gradient sign method (FGSM)** proposed in [1] generates adversarial examples by taking a single step:

$$X^{adv} = X + \epsilon * sign(\nabla_X L(X, y^{target}))$$

where the size of the perturbation  $\epsilon$  is often subject to some restrictions [7]. A common implementation of the FGSM attack is to gradually increase the magnitude of  $\epsilon$  until the input is misclassified. Its iterative extension named **basic iterative method** has the following update rule [18]:

$$X_0^{adv} = X, \qquad X_{N+1}^{adv} = Clip_{X,\epsilon} \left\{ X_N^{adv} + \alpha * sign(\nabla_X L(X_N^{adv}, y^{target})) \right\}$$

where  $\alpha$  regulates the size of the update on each step and the total size of perturbation is capped at  $\epsilon$  using  $Clip_{X,\epsilon}^{-1}$ .

## 2.2 White-box iterative attack example: DeepFool

The **DeepFool** attack is an iterative attack similar to the basic iterative method but also takes into account the  $\ell_2$ -norm of the gradients when computing the update rule. The original paper suggests that the proposed method generates adversarial perturbations which are hardly perceptible, while the fast gradient sign method outputs a perturbation image with higher norm [2]. Such patterns are observed on both the MNIST and CIFAR-10 dataset using the state-of-the-art architectures.

```
Algorithm 1: DeepFool Algorithm

Result: DeepFool Algorithm(binary case)

1 Initialize: x_0 \leftarrow x, i \leftarrow 0

2 while sign(f(x_i)) = sign(f(x_0)) do

3 \begin{vmatrix} r_i \leftarrow \frac{f(x_i)}{||\nabla f(x_i)||_2^2} \nabla f(x_i) \end{vmatrix}

4 x_{i+1} \leftarrow x_i + r_i

5 |i \leftarrow i + 1|

6 end
```

#### 2.3 Black-box iterative attack example: LocalSearch

In the case of black-box attacks, the adversaries do not have access to model architecture and have no internal knowledge of the target network. These kinds of methods treat the network as an oracle and

Here we borrow the notation from [7].  $Clip_{X,\epsilon}(A)$  clips A element-wise such that  $A_{i,j} \in [X_{i,j} - \epsilon, X_{i,j} + \epsilon]$ .

only assume that the output of the network can be observed on the probed inputs. The LocalSearch attack is accomplished by carefully constructing a small set of pixels to perturb by using the idea of greedy local search [3]. This is an extension of a simple adversarial attack, which randomly selects a single pixel and applies a strong perturbation to it in order to misclassify the input image. The LocalSearch attack is also an iterative procedure, where in each round a local neighborhood is used to refine the current image. This process minimizes the probability of assigning high confidence scores to the true class label, by the network. This approach identifies pixels with high saliency scores but without explicitly using any gradient information [3].

# 3 Defense: Methods Against Adversaries

On the defenders' side, proposed measures against adversarial attacks include input validation and preprocessing, adversarial training, defensive distillation and architecture modifications. In the paragraphs below, we briefly review these methods and discuss how randomized training/models such as stochastic delta rule could increase model robustness against adversaries.

## 3.1 Adversarial Training

Adversarial training increases model robustness by providing adversarial examples to the model during training. The standard practice is to generate adversarial examples from a subset of the incoming batch of clean inputs dynamically. The model is then trained on the mixed batch of clean and adversarial inputs. However in order to do this, a specific method for generating adversarial examples must be assumed, preventing adversarial training from being adaptive to different attack methods. For example, [7] showed that models adversarially trained using *one-step* methods are fooled easily by adversarial examples generated using *iterative* methods; models adversarially trained using a fixed  $\epsilon$  could even fail to generalize to adversarial examples created using different  $\epsilon$  values.

## 3.2 Defensive Distillation

Distillation was originally proposed in the context of model compression, aiming to transfer learned knowledge from larger, more complex models to more compact and computationally efficient models [8]. *Defensive distillation* was first proposed by [9] as a training regime to increase model robustness against adversaries. The goal of *defensive distillation* is not as much transfer learning (for which distillation was originally proposed), but rather to train models to have smoother gradient surfaces with respect to the input – such that small steps in the input space do not change the model's output significantly.

While smoothing out the gradients that adversaries usually use to create adversarial examples is effective in the setting described by the original paper, [19] pointed out that the attack assumed by [9] (Papernot's attack) could be oblivious to potentially stronger attacks. In addition, models trained with *defensive distillation* possess no advantage against a modified version of Papernot's attack when compared to regularly trained models. Works such as [20] investigate the effect of network compression solely for the purpose of transferring model knowledge, but discovers the effect of robustness against adversaries as a side product.

## 3.3 Randomized Methods & Models

Randomized training methods seek to improve the robustness of deep models by introducing randomness, irrespective of benign or adversarial samples, during the training process. For example, [12] introduces a random resizing layer and/or zero-padding layer prior to the regular architectures of CNNs. Through experimental evaluation, the authors discovered that this method is particularly effective against iterative attacks, while other methods introduced above are better at handling single-step attacks. A combination of both methods, as the authors argue, achieve best performance against arbitrary adversaries.

## 4 Adversarial robustness through stochastic parameters

Once trained, traditional neural networks have fixed parameters during inference. This permits the adversary to obtain consistent responses from the system, as well as consistent gradient values required to compute perturbations.

In our approach the network parameters are themselves stochastic, drawn from a distribution. Each parameter  $w_j$  in the network (which we assume without loss of generality to be a vector) has its own distribution  $P(w_j;\theta_j)$  with parameter  $\theta_j$ . When performing inference, the parameter value  $w_j$  is drawn from its distribution, i.e.  $w_j \sim P(w_j;\theta_j)$ . The process of training the network comprises learning the parameters of the distributions  $\theta_j$ , rather than the parameters  $w_j$  themselves.

As a consequence of the stochasticity of the network, the gradients computed by an adversary (for the purpose of generating adversarial samples) will actually be stochastic, and may not generalize to other runs inference when the drawn parameter values are different. While it may seem that this effect should have little influence and average out in expectation, particularly for black-box attacks, our experiments reveal that it is in fact sufficient to greatly decrease the efficacy of the adversary.

For our solution we use variants of the Stochastic Delta Rule [13] to build our network. We describe these below. We provide the specifics of the original SDR training routine (SDR-Decay), along with our proposed fine-grained variant of SDR (SDR-Learnable). Along this trajectory, we will raise the issue of practical implementation concerns, the connection between SDR and regularization, and a qualitative explanation of feasibility of a SDR-augmented training routine in improving model generalizability and robustness against adversaries.

#### 4.1 The Stochastic Delta Rule

First introduced in [13], SDR is revisited under the deep learning setting in [21]. During the forward pass of an SDR-equipped model, parameters  $w_j$  are not regarded as fixed values, but are rather random variables sampled from an arbitrary distribution  $P(w_j;\theta_j)$  specified by parameters  $\theta_j$ . The choice of such distribution is arbitrary. For the purpose of our experiments, we assume that model parameters  $w_j$  follow a normal distribution, i.e.  $P(w_j;\theta_j) = N(\mu_j, \Sigma_j)$ , and are independent of each other. The parameters  $\mu_j$  and  $\Sigma_j$  must be learned for each  $w_j$ .

In the discussions below on training these parameters, we will drop the subscript j for brevity. Given a (minibatch of) training input(s) (X,y) with features X and labels y, at training iteration t, a model with the SDR training routine samples model parameters  $w^{(t)} \sim N(\mu^{(t)}, \Sigma^{(t)})$ , fits the current batch with respect to the sampled parameters, and performs the following updates:

$$\mu^{(t+1)} \leftarrow \mu^{(t)} - \alpha \nabla_w L(X, y, w^{(t)})$$

$$\Sigma^{(t+1)} \leftarrow \Sigma^{(t)} + \beta |\nabla_w L(X, y, w^{(t)})|$$
(1)

where  $\alpha, \beta$  are step sizes for the mean and variance respectively, and, as before, L() is a loss function.

On one hand, it can be readily observed that with  $\beta=0$  and zero-intialization of the covariance matrix, only the mean parameters are updated, and we recover the usual training routine without SDR. On the other hand, [21] states that if parameters are sampled from a binomial distribution with mean Dp and variance Dp(1-p), SDR could be regarded as a special case of dropout with probability p, only in this case the variance is not updated with respect to information gathered from the gradients.

# 4.2 SDR Update with Scheduled Variance Decay

We first consider the SDR proposed by [21] which we formally present in algorithm 2. While updating model weights, the original SDR does not backpropagate gradients onto the parameter variances. Instead, the parameter variances are updated using the size of their associated means' gradients,  $|\nabla_w \partial L(X,y,w)|$  (see Equation 1). As explained by [21] in the original paper, the intuition behind this update rule is that a larger gradient on the parameter mean would motivate expansion of the random node's distribution to explore potentially better weight values. Besides this update by gradient, the original SDR anneals the parameter variances by  $\zeta$  to ensure asymptotically decaying variances (see Algorithm 2 below). As training progresses, variances are shrunk so as to sample progressively concentrated parameters around the mean, for which reason we term this original SDR as SDR-Decay. In our implementation, we further introduce the decay schedule,  $\tau$ , to avoid

over shrinkage and allow sufficient exploration of our model within the parameter space. At test time, under original the SDR, one would compute the forward pass directly with parameter means. However, in our case where SDR is applied to adversarial defense, doing so would defeat the purpose of generating stochastic gradients that mislead the adversary; thus we compute the forward pass of our model using weights sampled from the learned parameter distributions, as done during training.

Algorithm 2: SDR-Decay. The algorithm applies to every parameter in the network.

```
input dataset \{(X_i,y_i)\}_{i=1}^n, decay schedule \tau, decay rate \zeta\in(0,1], batch size B

Initialize model parameters \mu^{(0,0)}, \Sigma^{(0,0)}

num-batches \leftarrow n//B

for e=1,2,\cdots until convergence:

for b=1,2,\cdots, num-batches

Sample batch parameter weights w\sim N(\mu,\Sigma)

Compute forward pass of network with respect to w

Perform SDR parameter updates with respect to equations 1

If b\%\tau=0: \Sigma\leftarrow\zeta\Sigma

End for
```

#### 4.3 SDR-Learnable

Algorithm 2 updates both parameter means and variances with  $\nabla_w \partial L(X,y,w)$ , requiring backward pass on only one set of parameters w for each batch. We can, however, let both means and variances be fully learnable and have them updated with  $\nabla_w \partial L(X,y,w)$  and  $\nabla_\Sigma \partial L(X,y,w)$  respectively. The resulting algorithm, termed SDR-Learnable, produces a more realistic learning paradigm that approximates the behavior of variable parameters with higher fidelity. The drawback of such approach, of course, is doubling the computation required to compute the gradients. At test time, we sample parameters from learned means and variances, and perform inference on input data with sampled parameters.

It is worth noting that with this formulation of the update rule, we require the loss function to be differentiable with respect to the parameters  $\mu$  and  $\Sigma$ . However, the loss value is computed based on sampled realizations of many random variables along the forward pass. The sampling operation is not explicitly differentiable with respect to distribution parameters. In order to circumvent this issue, we use the reparametrization trick, as illustrated in figure 2, to make the network capable of backpropagating through random nodes. This technique essentially transfers the non-deterministic nature of the weight to another source of randomness, which then allows the randomly generated weights to be differentiable with respect to its parameters.

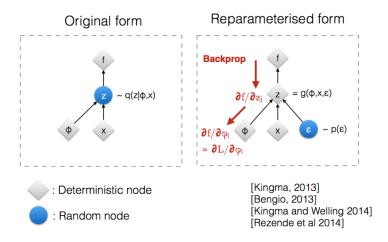


Figure 2: Reparameterization trick of SDR-Learnable update

In section 4.3.2, we show that the variance of the parameter distributions will shrink as the network is trained. Smaller variances result in less randomness in the forward pass, and as a result of this, the network does not exhibit consistent robustness to adversarial attacks. Furthermore, larger variances introduce a large amount of non-determinism in the forward pass. While this yields higher robustness, the overall task performance accuracy of the network diminishes.

In order to address this issue, we further modified the SDR architecture to use variance thresholds to ensure all parameter distributions' variances fall within a specified range during training. After performing each optimization step, we iterate through all parameters in the network and update them to fall within the desired range as shown in Algorithm 3 below.

## Algorithm 3: Variance Threshold

```
1 for \mu, \sigma in network parameters:

2 if \sigma > \max var:

3 \sigma = \max var

4 else if \sigma < \min var:

5 \sigma = \min var

6 End for
```

In order to find a trade-off between task-specific performance and adversarial robustness, we employed a training schedule to incrementally increase the variance of the parameter distributions. The minimum and maximum variances were initialized to 0.0 and 1.0 respectively. After every epoch of training, we evaluate the task-specific performance of the model on the test dataset. If the test performance begins to plateau, we increase both the minimum and maximum variances by 0.05. This update occurs at most once every five epochs of training. As we increase the variance, the test accuracy decreases initially. After fixing the variance thresholds, the test accuracy improves again on the subsequent training iterations. Increasing the minimum variance allows us to improve the non-determinsm, which directly impacts the network's robustness to adversarial examples. By controlling the maximum variance, we ensure that the network is still capable of achieving high task-specific performance. Using this training schedule, we are able to train a network with high variance that performs well on the clean data, and is robust to adversarial attack.

## 4.3.1 SDR-Learnable in Adversarial Learning

The inference procedure of SDR-Learnable produces variable predictions with the same input, which motivates us to investigate its robustness against adversarial samples. Our qualitative motivation for this hypothesis is as follows: adversarial attacks are designed to lead models to misclassify, while inducing no human-recognizable changes to inputs. The adversarially modified and original clean inputs must be very similar. To account for such perturbations, the decision boundaries of a model must not only cater to the specific data points provided in the training data, but also to a *vicinity* of these points in the sample space.

The variable treatment of parameters in SDR-Learnable is a step towards defensive strategies in two aspects: (1). the most effective gradient attack direction is computed with respect to one sampled parameter instance, and is less effective for another; (2) instead of fitting data at localized points, models with variable parameters create a decision region subject to the parameter distribution, hence allowing more robust prediction against adversarial samples.

#### 4.3.2 SDR-Learnable in Generalization

The description of decision regions, rather than spiky predictions, of SDR-Learnable naturally leads us to investigate its relationship to model generalization. In particular, we are interested in whether the inclusion of variable model parameters improves out-of-sample performance, and how do variances behave as training progresses. We present the following result as a first step towards the analysis.

Let  $X \in \mathbb{R}^{N \times D}$  and  $y \in \mathbb{R}^N$  be fixed, and  $w \in \mathbb{R}^D$  be a random vector with  $\mathbb{E}[w] = \mu$ ,  $\operatorname{Cov}[w] = \Sigma$ , then the risk of a linear regression model  $\hat{y} = Xw$  takes the form

$$\mathbb{E}_{w}[\|y - Xw\|^{2}] = \|y - X\mu\|^{2} + \|X\Sigma^{1/2}\|^{2}$$
(2)

With simple algebra of expectation, we observe that

$$\mathbb{E}_{w}[\|y - Xw\|^{2}] = \||y||^{2} - 2y^{T}X\mathbb{E}_{w}[w] + \sum_{i=1}^{N} \mathbb{E}_{w}[(X_{i}^{T}w)^{2}]$$
(3)

$$= ||y||^2 - 2y^T X \mathbb{E}_w[w] + \sum_{i=1}^N \left( \text{Var}[X_i^T w] + (X_i^T \mathbb{E}_w[w])^2 \right)$$
(4)

$$= ||y||^2 - 2y^T X \mu + \sum_{i=1}^{N} (X_i^2 \mu)^2 + \sum_{i=1}^{N} X_i^T \Sigma X_i$$
 (5)

$$= ||y - X\mu||^2 + ||X\Sigma^{1/2}||^2.$$
 (6)

It can be readily observed that under the stylized linear regression model, SDR-Learnable is equivalent to regularizing parameter variances  $\Sigma$  with penalty matrix X. In neural network training, we expect that the Frobenius norm of  $\Sigma$  to decay progressively, hence leading to more concentrated parameter samples. Notice that the decay step  $\Sigma \leftarrow \zeta \Sigma$  in SDR-Decay is a step towards artificial control of the magnitude of parameter variances, mimicking the behavior of SDR-Learnable.

It is worth noting that the parameter distribution does not necessarily lead to better generalization. Improved performance may be obtained by averaging the outcomes of multiple inferences, however this comes at the cost of adversarial robustness. Consequently we only perform a single pass of inference on any sample.

# 5 Experiments & Results

We experimentally evaluate the efficacy of the proposed SDR-Learnable in model generalization and robustness against adversarial samples on the MNIST [22] dataset, using the FoolBox toolkit [23] for adversarial samples. To compare performance, we evaluate three baseline models in addition to SDR-Learnable: a standard 3-layer MLP, a 3-layer SDR MLP, and a 3-layer MLP with DropConnect (p=0.2) [24]. Dropconnect is a generalization of dropout [25], which introduces non-determinism in the network by randomly dropping network connections. This non-determinism results in an improvement in adversarial robustness. The SDR-Learnable model used is a 3-layer MLP with learnable parameter distributions for weights and biases. All models had an input layer of size 784, followed by three hidden layers of size 100 and a final layer with 10 output neurons, representing the class probabilities. The ReLU activation was used on all layers except the output layer which used the Softmax activation. To train all models, we used the Cross-Entropy Loss function.

Experimental evaluation was restricted to these simple models as gradient-based white-box attacks are significantly harder to defend against in this setting; larger and more complex models' gradients w.r.t the input are more difficult to estimate, adversely affecting the generated adversarial sample. Furthermore, black-box attacks will be able to estimate a simple models' decision boundaries to a higher degree of accuracy, resulting in more challenging adversarial inputs.

The models were trained using the standard train and test split on the MNIST dataset[22], the test-set results are reported in Table.1. The classification accuracy on the regular test-set was consistent across all models, with SDR-Learnable achieving an overall accuracy of 97.87%, outperformining the MLP and SDR baselines. It shows superior robustness against adversarial samples. Particularly in the case of the one-step FGSM attack, the SDR-Learnable model achieves a classification accuracy that is comparable to the classification accuracy on the uncontaminated dataset.

	Vanilla MLP	SDR-decay MLP	Learnable SDR MLP	DropConnect
Regular Samples	97.59%	97.55%	97.87%	97.95%
FGSM Attack	0%	4.61%	94.86%	45.9%
DeepFool Attack	0%	3.87%	78.42%	45.48%
LocalSearch Attack	0%	2.85%	88.89%	26.53%

Table 1: Model performance on regular test set, and under adversarial attack.

## 6 Discussions and Conclusions

From the results in Table.1, we note that SDR-Learnable is robust against one-step and iterative white-box attacks. The decrease in classification accuracy between the FGSM and DeepFool attacks

is expected, as combating iterative attacks is a strictly harder task for adversarial defense[26]. Without using the variance threshold schedule to train SDR-Learnable, the classification accuracy under the DeepFool attack was 40.15%. This indicates that the use of the variance threshold technique is crucial to defense against iterative white-box attacks. SDR-Learnable achieves an adversarial accuracy of 88.89% against the iterative black-box attack LocalSearch. Without using the variance scheduler, the accuracy against LocalSearch is 44.05%. This indicates that there is significant benefit of using the variance scheduler to combat iterative black-box attacks.

While SDR-Learnable has been shown to be robust against adversarial attacks, it was not trained using any adversarial examples. Existing work has shown that the use of adversarial examples during training results in an adversarially robust network; the same technique can be leveraged using SDR-Learnable networks to further improve their robustness. Several existing techniques for adversarial defense can be incorporated into our network architecture, further improving model robustness.

In the reported experiments, SDR-Learnable was used only in an MLP, but the technique can be easily applied to almost any neural network architecture, resulting in at most twice the number of original parameters (each parameter is replaced with a distribution parameterized by a mean and a variance term). In our experiments, there was no significant difference in the time taken to train or evaluate the SDR-Learnable network.

In conclusion, we have demonstrated that non-determinism in the model parameters improves robustness against white-box and black-box attacks. The SDR-Learnable technique can be adapted using any network architecture, maintaining task-specific performance and providing defense against one-step and iterative whitebox and blackbox attacks, at no significant additional parameter cost. Furthermore, while iterative white-box attacks have been shown to compromise defense models based on stochastic gradients, we have shown that explicitly increasing the parameter variances while maintaining task-specific performance in SDR-Learnable significantly improves robustness.

## References

- [1] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and Harnessing Adversarial Examples. *arXiv e-prints*, page arXiv:1412.6572, Dec 2014.
- [2] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. *CoRR*, abs/1511.04599, 2015.
- [3] Nina Narodytska and Shiva Prasad Kasiviswanathan. Simple black-box adversarial perturbations for deep networks. *CoRR*, abs/1612.06299, 2016.
- [4] Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-gan: Protecting classifiers against adversarial attacks using generative models. *arXiv* preprint arXiv:1805.06605, 2018.
- [5] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083, 2017.
- [6] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defenses. *arXiv preprint arXiv:1705.07204*, 2017.
- [7] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. CoRR, abs/1611.01236, 2016.
- [8] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2015.
- [9] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. *CoRR*, abs/1511.04508, 2015.
- [10] Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *International Conference on Learning Representations*, 2018.
- [11] Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens van der Maaten. Countering adversarial images using input transformations. *arXiv preprint arXiv:1711.00117*, 2017.
- [12] et al. Xie, Cihang. Mitigating adversarial effects through randomization. arXiv: 1711. 01991, 2017.
- [13] Stephen José Hanson. "a stochastic version of the delta rule. *Physica D: Nonlinear Phenomena 42.1-3* (1990): 265-272., 1990.

- [14] Kevin Eykholt, Ivan Evtimov, Earlence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song. Robust physical-world attacks on deep learning models. arXiv preprint arXiv:1707.08945, 2017.
- [15] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference* on Computer and Communications Security, pages 1528–1540. ACM, 2016.
- [16] Samuel G Finlayson, Hyung Won Chung, Isaac S Kohane, and Andrew L Beam. Adversarial attacks against medical deep learning systems. *arXiv preprint arXiv:1804.05296*, 2018.
- [17] et al. Papernot, Nicolas. Practical black-box attacks against machine learning. *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, 2017.
- [18] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. CoRR, abs/1607.02533, 2016.
- [19] Nicholas Carlini and David A. Wagner. Defensive distillation is not robust to adversarial examples. CoRR, abs/1607.04311, 2016.
- [20] et al. Belagiannis, Vasileios. Adversarial network compression. arXiv: 1803. 10750, 2018.
- [21] Noah Frazier-Logue and Stephen José Hanson. Dropout is a special case of the stochastic delta rule: faster and more accurate deep learning. arXiv: 1808.03578, 2018.
- [22] Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [23] Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models. *arXiv preprint arXiv:1707.04131*, 2017.
- [24] Li Wan, Matthew Zeiler, Sixin Zhang, Yann Le Cun, and Rob Fergus. Regularization of neural networks using dropconnect. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th Inter*national Conference on Machine Learning, volume 28:3 of Proceedings of Machine Learning Research, pages 1058–1066, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- [25] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [26] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In Proceedings of the 35th International Conference on Machine Learning, ICML 2018, July 2018.