

# Unreasonable Effectiveness of Last Hidden Layer Activations for Adversarial Robustness

Omer Faruk Tuna

Computer Engineering Department  
Isik University  
Istanbul, Turkey  
omer.tuna@isikun.edu.tr

Ferhat Ozgur Catak

Electrical Engineering and Computer Science  
University of Stavanger  
Rogaland, Norway  
f.ozgur.catak@uis.no

M. Taner Eskil

Computer Engineering Department  
Isik University  
Istanbul, Turkey  
taner.eskil@isikun.edu.tr

**Abstract**—In standard Deep Neural Network (DNN) based classifiers, the general convention is to omit the activation function in the last (output) layer and directly apply the softmax function on the logits to get the probability scores of each class. In this type of architectures, the loss value of the classifier against any output class is directly proportional to the difference between the final probability score and the label value of the associated class. Standard White-box adversarial evasion attacks, whether targeted or untargeted, mainly try to exploit the gradient of the model loss function to craft adversarial samples and fool the model. In this study, we show both mathematically and experimentally that using some widely known activation functions in the output layer of the model with high temperature values has the effect of zeroing out the gradients for both targeted and untargeted attack cases, preventing attackers from exploiting the model’s loss function to craft adversarial samples. We’ve experimentally verified the efficacy of our approach on MNIST (Digit), CIFAR10 datasets. Detailed experiments confirmed that our approach substantially improves robustness against gradient-based targeted and untargeted attack threats. And, we showed that the increased non-linearity at the output layer has some additional benefits against some other attack methods like Deepfool attack.

**Index Terms**—trustworthy AI, adversarial machine learning, deep neural networks, robustness

## I. INTRODUCTION

By the end of 2013, researchers found out that DNN models are vulnerable to well-crafted malicious perturbations. Szegedy et al. [1] were the first to recognize the prevalence of adversarial cases in the context of image classification. Researchers have shown that a slight alteration in an image can influence the prediction of a DNN model. It is demonstrated that even the most advanced classifiers can be fooled by a very small and practically undetectable change in input, resulting in inaccurate classification. Since then, a lot of research studies [2]–[5] were performed in this new discipline known as *Adversarial Machine Learning* and these studies were not limited just to image classification task. To give some example, Sato et al. [6] showed in the NLP domain that changing just one word from an input sentence can fool a sentiment analyser trained with textual data. Another example is in the audio domain [7], where the authors generated targeted adversarial audio samples in autonomous speech recognition task by introducing very little distortion to the original waveform. The

findings of this study indicate that the target model can simply be exploited to transcribe the input as any desired phrase.

Attacks that take advantage of DNN’s weakness can substantially compromise the security of these machine learning (ML)-based systems, often with disastrous results. Adversarial evasion attacks mainly work by altering the input samples to increase the likelihood of making wrong predictions. These attacks can cause the model’s prediction performance to deteriorate since the model cannot correctly predict the actual output for the input instances. In the context of medical applications, a malicious attack could result in an inaccurate disease diagnosis. As a result, it has the potential to impact the patient’s health, as well as the healthcare industry [8]. Similarly, self-driving cars employ ML to navigate traffic without the need for human involvement. A wrong decision for the autonomous vehicle based on an adversarial attack could result in a tragic accident. [9], [10]. Hence, defending against malicious evasion attacks and boosting the robustness of ML models without sacrificing clean accuracy is critical. Presuming that these ML models are to be utilized in crucial areas, we should pay utmost attention to ML models’ performance and the security problems of these architectures.

In principle, adversarial strategies in evasion attacks can be classified based on multiple criteria. Based on the attacker’s ultimate goal, attacks can be classified as targeted and untargeted attacks. In the former, the attacker perturbs the input image, causing the model to predict a class other than the actual class. Whereas in the latter, the attacker perturbs the input image so that a particular target class is predicted by the model. Attacks can also be grouped based on the level of knowledge that the attacker has. If the attacker has complete knowledge of the model like architecture, weights, hyper-parameters etc., we call this kind of setting as White-Box Settings. However, if the attacker has no information of the deployed model and defense strategy, we call this kind of setting as Black-Box Settings [11].

This research study focused on both targeted and untargeted attacks in a White-Box setting. We propose an effective modification to the standard DNN-based classifiers by adding special kind of non-linear activation functions (sigmoid or tanh) to the last layer of the model architecture. We showed that training a model using high temperature value at the

output layer activations and using the model by discarding the temperature value at inference time provides a very high degree of robustness to loss-based White-box targeted and untargeted attacks, together with attacks acting like Deepfool. We hereby name our proposed models as Squeezed Models. Our codes are released on GitHub<sup>1</sup> for scientific use.

To summarize; our main contributions for this study are:

- We propose an effective modification to standard DNN based classifiers, which enables natural robustness to gradient-based White-box targeted and untargeted attacks.
- We show that using a specific type of non-linear activation functions at the output layer with high temperature values can actually provide robustness to the model without impairing the ability to learn.
- We experimentally showed that adding non-linearity to the last hidden layer provides robustness to other types of attacks, like Deepfool.

## II. RELATED WORKS

Since the uncovering of DNN’s vulnerability to adversarial attacks [1], a lot of work has gone into inventing new adversarial attack algorithms and defending against them by utilizing more robust architectures [12]–[15]. We will discuss some of the noteworthy attack and defense studies separately.

### A. Adversarial evasion attacks

DNN models have some vulnerabilities that make them challenging to defend in adversarial settings. For example, they are mostly sensitive to slight changes in the input data, leading to unexpected results in the model’s predictions. Figure 1 depicts how an adversary could take advantage of such a vulnerability and fool the model using properly crafted perturbation applied to the input.

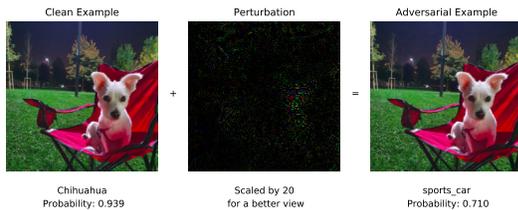


Fig. 1: The figure depicts an example to adversarial attack. The original image is subjected to the adversarial perturbation. The precisely crafted perturbation manipulates the model in such a way that a ”Dog (Chihuahua)” is wrongly identified as ”Sports Car” with high confidence.

An important portion of the attack methods are gradient-based and based on perturbing the input sample in order to

maximize the model’s loss. In recent years, many different adversarial attack techniques have been suggested in literature. The most widely known and used adversarial attacks are Fast-Gradient Sign Method, Iterative Gradient Sign Method, DeepFool and Carlini-Wagner. These adversarial attack algorithms are briefly explained in Section II-A1 - II-A4.

1) *Fast-Gradient Sign Method*: This method, also known as FGSM [16], was one of the first and most well-known adversarial attacks. The derivative of the model’s loss function with respect to the input sample is exploited in this attack strategy to determine which direction the pixel values in the input image should be changed in order to minimize the loss function of the model. Once this direction is determined, it changes all pixels in the opposite direction at the same time to maximize loss. One can craft adversarial samples for a model with a classification loss function represented as  $J(\theta, \mathbf{x}, y)$  by utilizing the formula below, where  $\theta$  denotes the parameters of the model,  $\mathbf{x}$  is the benign input, and  $y_{true}$  is the real label of our input.

$$\mathbf{x}^{adv} = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y_{true})) \quad (1)$$

In [17], the authors presented a targeted variant of FGSM referred to as the Targeted Gradient Sign Method (TGSM). This way, they could change the attack to try to convert the model’s prediction to a particular class. To achieve this, instead of maximizing the loss with respect to the true class label, TGSM attempts to minimize the loss with respect to the target class  $J_{target}$ .

$$\mathbf{x}^{adv} = \mathbf{x} - \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y_{target})) \quad (2)$$

Different from Eq. 1, we now subtract the crafted perturbation from the original image as we try to minimize the loss this time. If we want to increase the efficiency of this approach, we can modify above equation as in Eq.3. The only difference is that instead of minimizing the loss of the target label, we maximize the loss of the loss of the true label and also minimize the loss for the alternative label.

$$\mathbf{x}^{adv} = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} (J(\theta, \mathbf{x}, y_{true}) - J(\theta, \mathbf{x}, y_{target}))) \quad (3)$$

2) *Iterative Gradient Sign Method*: Kurakin et al. [17] proposed a minor but significant enhancement to the FGSM. Instead of taking one large step  $\epsilon$  in the direction of the gradient sign, we take numerous smaller steps  $\alpha$  and utilize the supplied value  $\epsilon$  to clip the output in this method. This method is also known as the Basic Iterative Method (BIM), and it is simply FGSM applied to an input sample iteratively. Equation 4 describes how to generate perturbed images under the  $l_{inf}$  norm for a BIM attack.

$$\begin{aligned} \mathbf{x}_t^* &= \mathbf{x} \\ \mathbf{x}_{t+1}^* &= \text{clip}_{x,\epsilon} \{ \mathbf{x}_t + \alpha \cdot \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}_t^*, y_{true})) \} \end{aligned} \quad (4)$$

<sup>1</sup><https://github.com/author-name/xxx>

where  $\mathbf{x}$  is the clean sample input to the model,  $\mathbf{x}^*$  is the output adversarial sample at  $i^{\text{th}}$  iteration,  $J$  is the loss function of the model,  $\theta$  denotes model parameters,  $y_{true}$  is the true label for the input,  $\epsilon$  is a configurable parameter that limits maximum perturbation amount in given  $l_{inf}$  norm, and  $\alpha$  is the step size.

As in the case of TGSM, we can easily modify Eq. 4 to produce targeted variant of BIM. At each intermediate step, we can try to minimize the loss with respect to target class while at the same time maximizing the loss with respect to original class as in Eq. 5.

$$\begin{aligned} \mathbf{x}_t^* &= \mathbf{x} \\ \text{arxiv}\mathbf{x}_{t+1}^* &= \text{clip}_{x,\epsilon}\{\mathbf{x}_t + \alpha \cdot \text{sign}(\nabla_{\mathbf{x}}(J(\theta, \mathbf{x}_t^*, y_{true}) - \\ &\quad J(\theta, \mathbf{x}_t^*, y_{target})))\} \end{aligned} \quad (5)$$

3) *Deepfool Attack*: This attack method has been introduced by Moosavi-Dezfooli et al. [18] and it is one of the most strong untargeted attack algorithms in literature. It's made to work with several distance norm metrics, including  $l_{inf}$  and  $l_2$  norms.

The Deepfool attack is formulated on the idea that neural network models act like linear classifiers with classes separated by a hyperplane. Starting with the initial input point  $\mathbf{x}_t$ , the algorithm determines the closest hyperplane and the smallest perturbation amount, which is the orthogonal projection to the hyperplane, at each iteration. The algorithm then computes  $\mathbf{x}_{t+1}$  by adding the smallest perturbation to the  $\mathbf{x}_t$  and checks for misclassification. The illustration of this attack algorithm is provided in Figure 2. This attack can break defensive distillation method and achieves higher success rates than previously mentioned iterative attack approaches. But the downside is, produced adversarial sample generally lies close to the decision boundary of the model.

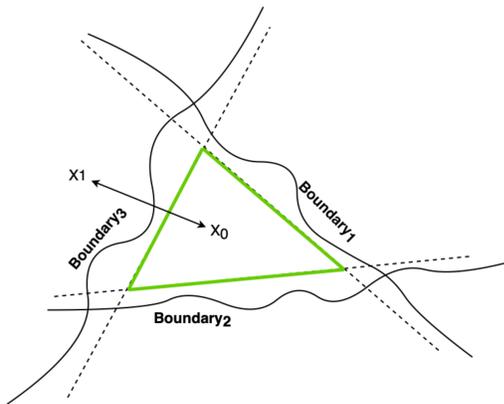


Fig. 2: Illustration of Deepfool attack algorithm

4) *Carlini&Wagner Attack*: Proposed by Carlini and Wagner [19], and it is one of the strongest attack algorithms so far. As a result, it's commonly used as a benchmark for the adversarial defense research groups, which tries to develop more robust DNN architectures that can withstand adversarial attacks. It is shown that, for the most well-known datasets,

the CW attack has a greater success rate than the other attack types on normally trained models. Like Deepfool, it can also deceive defensively distilled models, which other attack types struggle to create adversarial examples for.

In order to generate more effective and strong adversarial samples under multiple  $l_p$  norms, the authors reformulate the attack as an optimization problem which may be solved using gradient descent. A *confidence* parameter in the algorithm can be used to change the level of prediction score for the created adversarial sample. For a normally trained model, application of CW attack with default setting (confidence set to 0) would generally yield to adversarial samples close to decision boundary. And high confident adversaries generally located further away from decision boundary.

Adversarial machine learning is a burgeoning field of research, and we see a lot of new adversarial attack algorithms being proposed. Some of the recent remarkable ones are: i) Square Attack [20] which is a query efficient black-box attack that is not based on model's gradient and can break defenses that utilize gradient masking, ii) HopSkipJumpAttack [21] which is a decision-based attack algorithm based on an estimation of model's gradient direction and binary-search procedure for approaching the decision boundary, iii) Prior Convictions [3] which utilizes two kinds of gradient estimation (time and data dependent priors) and propose a bandit optimization based framework for adversarial sample generation under loss-only access black-box setting and iv) Uncertainty-Based Attack [2] which utilizes both the model's loss function and quantified epistemic uncertainty to generate more powerful attacks.

## B. Adversarial defense

1) *Defensive Distillation*: Although the idea of knowledge distillation was previously introduced by Hinton et al. [22] to compress a large model into a smaller one, the utilization of this technique for adversarial defense purposes was first suggested by Papernot et al. [23]. The algorithm starts with training a *teacher model* on training data by employing a high temperature ( $T$ ) value in the softmax function as in Equation 6, where  $p_i$  is the probability of  $i^{\text{th}}$  class and  $z_i$ 's are the logits.

$$p_i = \frac{\exp(\frac{z_i}{T})}{\sum_j \exp(\frac{z_j}{T})} \quad (6)$$

Then, using the previously trained teacher model, each of the samples in the training data is labeled with soft labels calculated with temperature ( $T$ ) in prediction time. The *distilled model* is then trained with the soft labels acquired from the teacher model, again with a high temperature ( $T$ ) value in the softmax. When the training of the student model is over, we use temperature value as 1 during prediction time. Figure 3 shows the overall steps for this technique.

2) *Adversarial Training*: Adversarial training is an intuitive defense method in which the model's robustness is increased by training it with adversarial samples. As demonstrated in Eq. 7, this strategy can be mathematically expressed as a Minimax game.

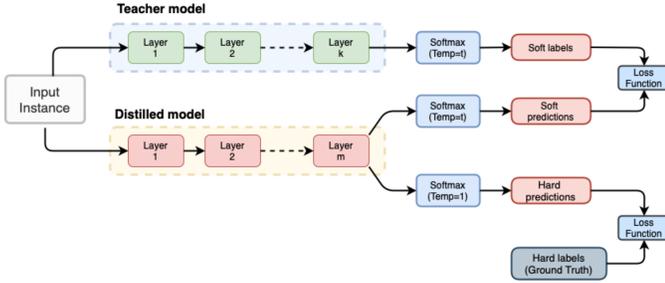


Fig. 3: Defensive Distillation.

$$\min_{\theta} \max_{|\delta| \leq \epsilon} J(h_{\theta}(x + \delta), y) \quad (7)$$

where  $h$  denotes the model,  $J$  denotes the model's loss function,  $\theta$  represents model's weights and  $y$  is the actual label.  $\delta$  is the amount of perturbation amount added to input  $x$  and it is constrained by given  $\epsilon$  value. The inner objective is maximized by using the most powerful attack possible, which is mostly approximated by various adversarial attack types. In order to reduce the loss resulting from the inner maximization step, the outside minimization objective is used to train the model. This whole process produces a model that is expected to be resistant to adversarial attacks used during the training of the model. For adversarial training, Goodfellow et al. [16] used adversarial samples crafted by the FGSM attack. And Madry et al. used the PGD attack [24] to build more robust models, but at the expense of consuming more computational resources. Despite the fact that adversarial training is often regarded as one of the most effective defenses against adversarial attacks, adversarially trained models are nevertheless vulnerable to attacks like CW.

Adversarial ML is a very active field of research, and new adversarial defense approaches are constantly being presented. Among the most notable are: i) High-Level Representation Guided Denoiser (HGD) [25] which avoids the error amplification effect of a traditional denoiser by utilizing the error in the upper layers of a DNN model as loss function and manages the training of a more efficient image denoiser, ii) APE-GAN [26] which uses a Generative Adversarial Network (GAN) trained with adversarial samples to eliminate any adversarial perturbation of an input image, iii) Certified Defense [27] which proposes a new differentiable upper bound yielding a model certificate ensuring that no attack can cause the error to exceed a specific value and iv) [4] which uses several uncertainty metrics for detecting adversarial samples.

### III. APPROACH

#### A. Chosen Activation Functions

We used specific type of activation functions (sigmoid and hyperbolic tangent) whose derivatives can be expressed in terms of themselves. And the derivative of these activation functions approaches to 0 when the output of the activation functions approaches to their maximum and minimum values.

Starting with the sigmoid function; we know that sigmoid function ( $\sigma(x)$ ) can be represented as in Eq. 8 and it squeezes the input to the range of 0 and 1 as can be seen in Figure 4a.

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (8)$$

And the derivative of Sigmoid function can be expressed as in Eq. 9:

$$\frac{d}{dx} \sigma(x) = \sigma(x) \cdot (1 - \sigma(x)) \quad (9)$$

One can easily derive using above formulation or verify from Figure 4a that the derivative of sigmoid function approach to 0 when the output of the sigmoid function approaches to 0 or 1.

Similarly, we can represent hyperbolic tangent ( $\tanh(x)$ ) function as in Eq. 10. Different from sigmoid, hyperbolic tangent function squeezes the input to the range of -1 and 1 as can be seen in Figure 4b.

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (10)$$

The derivative of hyperbolic tangent function can be expressed as in Eq. 11. Using Eq. 11 or Figure 4b, we can verify that the derivative of tanh function approaches to 0 when the output of the tanh function approaches to -1 or 1. So, the pattern is similar to the one we see in sigmoid function. The derivative of both of these activation functions yields to 0 when the output of the activation functions are at their minimum or maximum values. This property will be quite usefull when use these activation functions at the output layer of DNN classifiers to zeroing out the gradients.

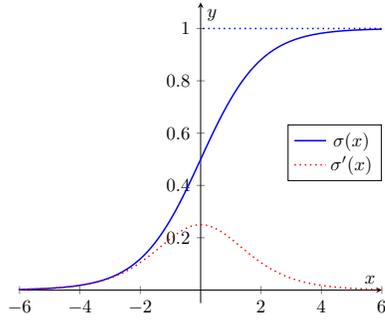
$$\frac{d}{dx} \tanh x = 1 - \tanh^2 x \quad (11)$$

#### B. Proposed Method

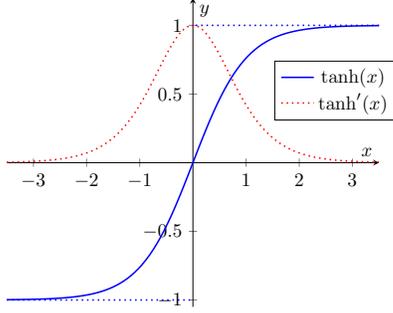
We begin this part by introducing the loss calculation for a standard deep neural network classifier. Let  $K$  denotes the number of output classes,  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$  be our dataset, where  $x_i \in \mathbb{R}^d$  and  $y_i \in \{o_1, o_2, \dots, o_k\}$  are the  $i^{th}$  input and output respectively where  $o_k$  is the one-hot encoded vector with the only  $k^{th}$  index being one and zero for the other indices and the probability output score of any output class with index  $k \in \{0, 1, \dots, K-1\}$  is represented by  $P_k$ . Based on this notation, the loss value ( $J$ ) of the classifier for any test input  $x^*$  can be calculated using cross-entropy loss function as below:

$$J = - \sum_{k=0}^{K-1} o_k[k] \cdot \log(P_k) = -\log(P_{true}) \quad (12)$$

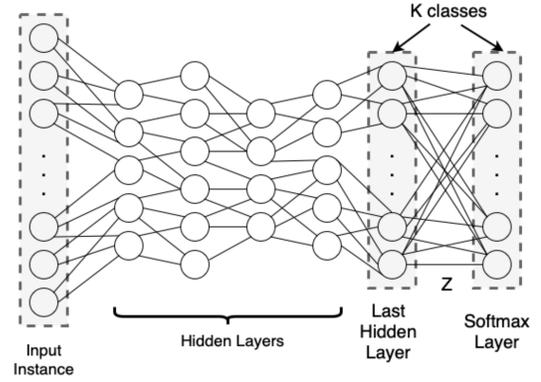
As can be seen in Figure 5a, in standard DNN-based classifiers that are widely used today, usage of activation functions in the output layer is omitted and the prediction scores of each class is calculated by feeding the output of the



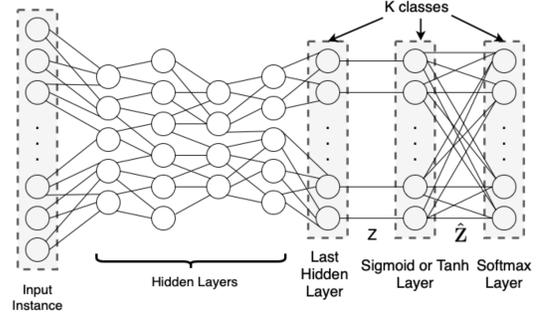
(a) Sigmoid



(b) Hyperbolic Tangent



(a) Standard DNN Classifier



(b) The proposed classifier

last layer of the network (logits) to the softmax function. If we denote the logits by  $Z = \{z_0, z_1, \dots, z_{K-1}\}$ , we can calculate the derivative of the loss with respect to  $k^{th}$  logit using Eq. 13. Formal derivation of the Eq. 13 is provided in Appendix B.

$$\frac{\partial J}{\partial z_k} = P_k - o_k[k] \quad (13)$$

Loss-based adversarial attacks try to exploit the gradient of the loss function ( $J$ ) with respect to input sample  $x$ , and what the attacker is trying to do is to use  $\frac{\partial J}{\partial x}$  to maximize  $J$ . We know from chain rule that  $\frac{\partial J}{\partial x} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial x}$ . Therefore, for any target class  $k$ , the gradient of the model's loss function with respect to the input image is directly proportional to  $\frac{\partial J}{\partial z_k}$ . In response to such kind of attack idea, several defense approaches have been proposed which mask the gradients of the models. For example, Defensive Distillation technique achieves this against untargeted loss-based attacks by enabling the model to make highly confident predictions. Because, when the model makes highly confident predictions in favor of the true class;  $P_{true}$  approaches 1, and since the label for true class is also 1,  $\frac{\partial J}{\partial z}$  and therefore  $\frac{\partial J}{\partial x}$  approaches to 0 for the specific untargeted attack case.

However, above approach will not work for targeted attack case. Because, in order to prevent targeted attacks, we have to make  $\frac{\partial J}{\partial z}$  to become 0 for target class. And, the way to achieve this for standard DNN-based classifiers is to make target probability ( $P_{target}$ ) to be very close to 1 (to make  $P_{target} - o_{target}[target]$  equals to 0) which obviously contradicts with the natural learning task. Therefore, there actually exists a dilemma between masking the gradient of the model

Fig. 5: Comparison of standard DNN classifier and the proposed classifier

for targeted attack case and achieving the task of learning for the model at our hand. This phenomenon is beautifully explained by Katzir et al. in [28].

To overcome this problem, we propose to use either of the two commonly known nonlinear activation functions (sigmoid and tanh) on the logits of the model as depicted in Figure 5b. The important thing is to apply an high temperature value to these activation functions during learning process (e.g. :  $\sigma(x, T) = 1/(1 + \exp(-x/T))$ ) and use the model by ignoring the temperature value at prediction time, just like defensive distillation technique. After our proposed modification, the output of the last layer will be  $\hat{Z}$ , where  $\hat{Z} = \{\hat{z}_0, \hat{z}_1, \dots, \hat{z}_{K-1}\}$  and  $\hat{Z} = \tanh(Z)$  or  $\hat{Z} = \sigma(Z)$ , depending on the chosen activation function. Based on this modified architecture, derivative of the model's loss with respect to input image under gradient-based attack against any class  $k$  can be formulated as below:

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial \hat{z}_k} \cdot \frac{\partial \hat{z}_k}{\partial z_k} \cdot \frac{\partial z_k}{\partial x} \quad (14)$$

In case of sigmoid function, the above equation can be reformulated as below by using Eq. 9 and Eq.13.

$$\frac{\partial J}{\partial x} = (P_k - o_k[k]) \cdot \hat{z}_k \cdot (1 - \hat{z}_k) \cdot \frac{\partial z_k}{\partial x} \quad (15)$$

And in case of tanh function, the Eq.14 can be written as below Eq. 11 and Eq.13:

$$\frac{\partial J}{\partial x} = (P_k - o_k[k]) \cdot (1 - \hat{z}_k^2) \cdot \frac{\partial z_k}{\partial x} \quad (16)$$

During the training of the DNN classifier depicted in Figure 5b, we force  $\hat{z}_k$  to be at it's maximum possible value for the true class in order to maximize the final softmax prediction score. And similarly, we force  $\hat{z}_k$  to be at it's minimum value for the other classes. Therefore, in case of sigmoid and tanh functions,  $\hat{z}_k$  will approach to 1 for true class. And for the classes other than the true class,  $\hat{z}_k$  will approach to 0 and -1 for sigmoid and tanh functions respectively. Since we additionally applied a high temperature value to these activation functions during training time, the output of these activation functions ( $\hat{z}_k$ ) will be even more close to their maximum and minimum values at prediction time when we omit their temperature values. Consequently, Eq. 15 and Eq. 16 will approach to 0 for both targeted and untargeted attack cases. Because, if we use the proposed model architecture using sigmoid function,  $\hat{z}_k \cdot (1 - \hat{z}_k)$  will be 0 when  $\hat{z}_k$  is either 0 or 1. And if we use the proposed model architecture using tanh function,  $1 - \hat{z}_k^2$  will become 0 when  $\hat{z}_k$  is either -1 or 1. This way, we can successfully zero out (mask) the gradients of the model for loss-based targeted and untargeted attack threats. To avoid any round-off errors in floating point operations, high precision should be set for floating point numbers in the ML calculations.

### C. Visual Representations of Loss Surfaces

We know that normally-trained models are vulnerable to gradient-based white-box targeted and untargeted attack threats. The main reason for this vulnerability lies in the ability of the attacker to successfully exploit the loss function of the model. To illustrate this fact, we made a simple experiment using a test image from MNIST (Digit) dataset and draw the loss surfaces of various models against two different directions (one for loss gradient direction and one for a random direction). When we check Figures 6a and 6b which display the loss surfaces of normally-trained model, we see in both cases that there exists a useful gradient information which the attacker can exploit to craft both untargeted and targeted adversarial samples.

To prevent the above vulnerability, various defense methods have been proposed, including Defensive distillation. This technique was found to significantly reduce the ability of traditional gradient-based untargeted attacks to build adversarial samples. Because defensive distillation has an effect of diminishing the gradients down to zero for untargeted attack case and the usage of standard objective function is not effective anymore. As depicted in Figure 6c, the gradient of the distilled model diminishes to zero and thus loss-based untargeted attacks have difficulty in crafting adversarial samples for defensively distilled models. However, it was then demonstrated that attacks, such as the TGSM attack,

could defeat the defensive distillation strategy [29], but without providing a mathematical proof about why these attacks actually work. And the actual reason of success for these kind of attacks against defensively distilled model is shown to lie in the targeted nature of these attacks [28]. Figure 6d demonstrates the loss surface of a distilled model under a targeted attack and we can easily see that the gradient of the model loss does not diminish to zero as in Figure 6c. The result is not surprising at all, because for a defensively distilled model under targeted attack, we expect  $P_{target}$  to be almost 0 and  $o_{target}[target]$  is 1. Therefore, we expect  $\frac{\partial J}{\partial z}$  (equivalent to  $P_{target} - o_{target}[target]$ ) approach to -1 which is more than enough to exploit the gradient of the loss function for a successful attack.

As a last attempt, we analyze the loss surfaces of one of our proposed models (model which was trained using tanh activation function with high temperature value at the output layer). When we check Figures 6e and 6f, we see that the gradient of model's loss function diminishes to 0 for both of the untargeted and targeted attack cases. And this prevents the attacker to exploit the gradient information of the model to craft successful adversarial perturbations.

### D. Softmax prediction scores of proposed architectures

For a normally trained standard DNN-based classifiers, we expect the model to make a prediction in favor of true class with a prediction score usually close to 1. In case of a defensively distilled model, we force the model to make high confident predictions. That is why, we see a prediction score very close to 1 in favor of the true class. However, for our proposed model architectures, the softmax prediction score of the true class is lower compared to a normal or defensively distilled model. Because the activation functions in the last layer of the model limits the values for  $\hat{z}_k$  to an interval of (0,1) and (-1,1). If we use sigmoid function in the last layer, maximum prediction score will be 0.232 and if we use tanh function in the last layer, maximum prediction score will be 0.450. And this will be the case for all the predictions. Similarly, minimum prediction scores will be 0.085 and 0.061 for models with sigmoid and tanh activation functions respectively. Softmax prediction score output of a test sample from MNIST dataset is displayed in Figure 7 for various models. We believe that this behaviour of our models, just like the case in defensively distilled models, might also be quite useful to prevent attackers to infer an information that is supposed to be private from the output probability scores of any prediction of the model and might contribute to the privacy of model as suggested in [30], [31].

## IV. EXPERIMENTS

### A. Adversarial Assumptions

In this research study, we assume that the attacker can choose to implement targeted or untargeted attacks towards the target model. Our assumption was that the attacker was fully aware of the architecture and parameters of the target model as in the case of *whitebox* setting and use the model as it

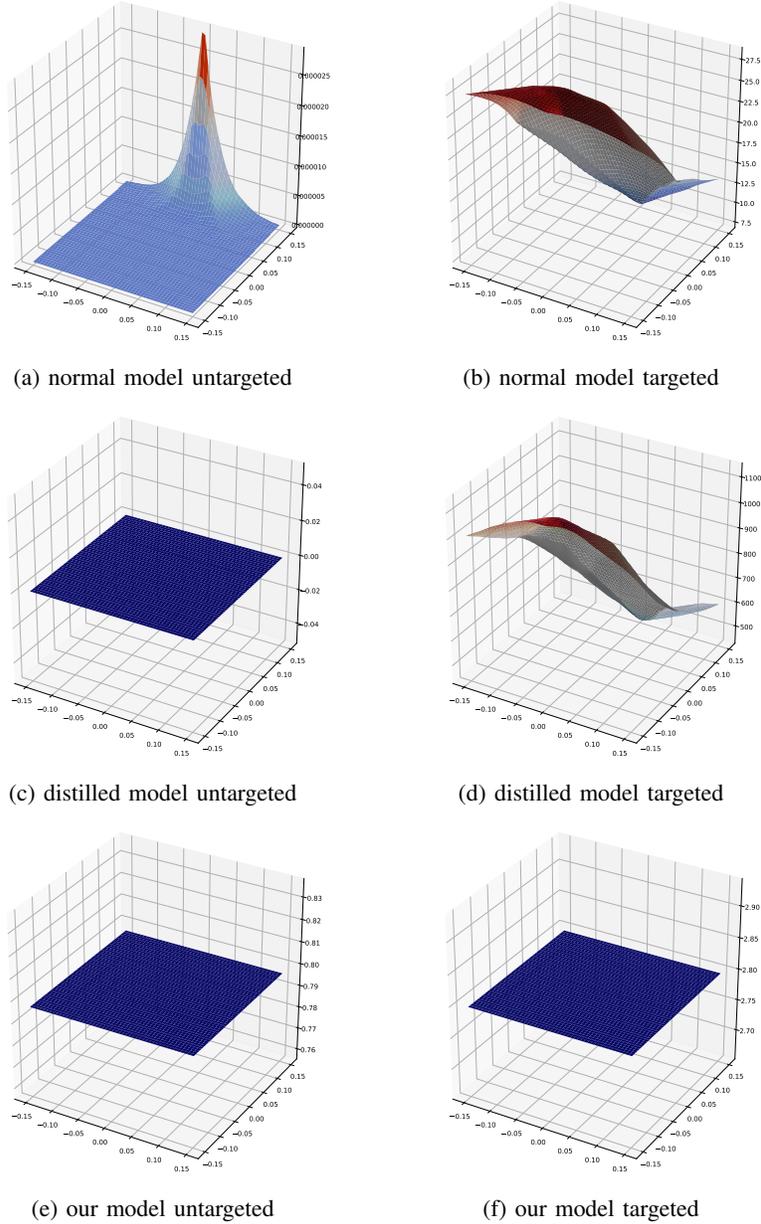


Fig. 6: Loss surfaces of various models under untargeted and targeted attack scenarios

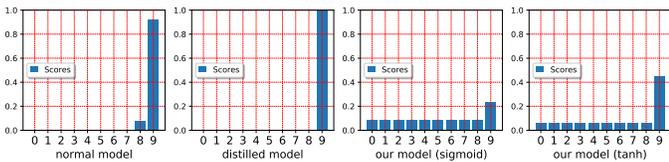


Fig. 7: Softmax score outputs of various models

is. Another crucial assumption concerns the constraints of the attacker. Clearly, the attacker should be limited to applying a perturbation with  $l_p$  norm up to certain  $\epsilon$  value for an attack to be unrecognizable to the human eye. For this study, we used  $l_\infty$  and  $l_2$  norm metrics to restrict the maximum perturbation

amount that an adversary can apply on the input sample. Finally, the error rate of our proposed defense technique is assessed over the percentage of resulting successful attack samples which is proposed by Goodfellow et al. [16] and recommended by Carlini et al. [32].

### B. Experimental Setup

For our experiments, we used four sets of models for each dataset as normal model, defensively distilled(student) model, proposed model with sigmoid activation and proposed model with tanh activation function at the output layer. By using same architectures, we trained our CNN models using MNIST (Digit) [33] and CIFAR-10 [34] datasets. For MNIST

TABLE I: Attack success rates on MNIST (Digit) - Part 1

	Normal Model	Distilled Model	Our Model (Sigmoid)	Our Model (Tanh)
FGSM ( $l_\infty, \epsilon : 0.1$ )	9.75%	2.13%	0.12%	0.38%
TGSM ( $l_\infty, \epsilon : 0.1$ )	1.77%	1.74%	0.04%	0.03%
BIM ( $l_\infty, \epsilon : 0.1$ )	34.20%	2.31%	0.05%	0.19%
Targeted-BIM ( $l_\infty, \epsilon : 0.1$ )	13.04%	9.31%	0.03%	0.02%
CW ( $l_2, \epsilon : 1.35 \text{ conf} : 0$ )	80.94%	59.99%	0.04%	0.07%
Deepfool ( $l_2, \epsilon : 1.35$ )	29.73%	21.22%	0.06%	0.14%

TABLE II: Attack success rates on MNIST (Digit) - Part 2

	Normal Model	Distilled Model	Our Model (Sigmoid)	Our Model (Tanh)
FGSM ( $l_\infty, \epsilon : 0.2$ )	31.09%	2.23%	0.38%	0.12%
TGSM ( $l_\infty, \epsilon : 0.2$ )	9.86%	8.33%	0.03%	0.04%
BIM ( $l_\infty, \epsilon : 0.2$ )	98.19%	2.71%	0.23%	0.08%
Targeted-BIM ( $l_\infty, \epsilon : 0.2$ )	90.05%	77.77%	0.03%	0.04%
CW ( $l_2, \epsilon : 2.70 \text{ conf} : 0$ )	100%	99.96%	0.11%	0.11%
Deepfool ( $l_2, \epsilon : 2.70$ )	97.69%	87.41%	0.17%	0.06%

(Digit) dataset, our models attained accuracy rates of 99.35%, 99.41%, 98.97% and 99.16% respectively. And, for CIFAR-10 dataset, our models attained accuracy rates of 83.95%, 84.68%, 82.37% and 80.15% respectively. The architectures of our CNN models and the hyperparameters used in model training are listed Appendix A. Finally, we set the temperature ( $T$ ) value as 20 and 50 for MNIST and CIFAR datasets respectively during the training of the defensively distilled model and our proposed models.

### C. Experimental Results

During our tests, we only implemented attack on the test samples if our models had previously classified them accurately. Because an adversary would have no reason to tamper with samples that have already been labeled incorrectly. For the TGSM and Targeted BIM attacks, we regard the attacks successful only if each the perturbed image is classified by the model as the chosen target class. We set the target class to "2" for MNIST (Digit) dataset, and "Cars" for CIFAR-10 dataset. We utilized an open source Python library called Foolbox [35] to implement the attacks used in this study. The attack parameters used in BIM and Targeted BIM are provided in Table VII.

The results of our experiments for MNIST and CIFAR10 datasets are available in Tables I,II and Tables III,IV together with the amount of perturbations applied and chosen norm metrics. Just for CW and Deepfool attacks, we used the  $l_2$  norm equivalent of the applied perturbation by using the formula  $l_2 = l_{inf} \times \sqrt{n} \times \sqrt{2}/\sqrt{\pi e}$ , where  $n$  is the input sample dimension. When we check the results, we observe that normally trained models are vulnerable to both targeted and untargeted attack types, whereas defensively distilled models are vulnerable to only targeted attack types. And our proposed models (squeezed models) provides a high degree of robustness to both targeted (TGSM, Targeted BIM, CW) and untargeted (FGSM, BIM) attacks. This success results from the effectiveness of our models in zeroing out the gradients in both scenarios.

TABLE III: Attack success rates on CIFAR10 - Part 1

	Normal Model	Distilled Model	Our Model (Sigmoid)	Our Model (Tanh)
FGSM ( $l_\infty, \epsilon : 3/255$ )	72.38%	13.88%	4.07%	1.64%
TGSM ( $l_\infty, \epsilon : 3/255$ )	22.84%	21.36%	0.56%	0.23%
BIM ( $l_\infty, \epsilon : 3/255$ )	93.53%	15.01%	2.61%	0.95%
Targeted-BIM ( $l_\infty, \epsilon : 3/255$ )	57.36%	57.22%	0.46%	0.22%
CW ( $l_2, \epsilon : 0.798$ )	100.00%	100.00%	2.93%	1.36%
Deepfool ( $l_2, \epsilon : 0.798$ )	99.98%	99.76%	2.22%	0.87%

TABLE IV: Attack success rates on CIFAR10 - Part 2

	Normal Model	Distilled Model	Our Model (Sigmoid)	Our Model (Tanh)
FGSM ( $l_\infty, \epsilon : 6/255$ )	80.58%	13.85%	4.04%	1.7%
TGSM ( $l_\infty, \epsilon : 6/255$ )	25.41%	25.58%	0.64%	0.3%
BIM ( $l_\infty, \epsilon : 6/255$ )	96.75%	15.11%	3.07%	1.14%
Targeted-BIM ( $l_\infty, \epsilon : 6/255$ )	68.16%	73.07%	0.6%	0.29%
CW ( $l_2, \epsilon : 1.596$ )	100%	100%	3.08%	1.75%
Deepfool ( $l_2, \epsilon : 1.596$ )	99.98%	100%	2.17%	0.89%

One other thing worth to mention about the result of our experiments is that, in addition to gradient-based attacks, our proposed models exhibit excellent performance against Deepfool attack as well. Generally, the reason behind the success of Deepfool attack against standard DNN-based classifiers is the linear nature of these models as argued by Goodfellow et. al. [16] and the authors of Deepfool paper formalized their methods based on this assumption [18]. However, since we introduce additional non-linearity to the standard DNN classifiers at the output layer, Deepfool attack algorithm fails to succeed in crafting adversarial samples compared to normally or defensively distilled models.

### V. CONCLUSION

In this study, we first showed that existing DNN-based classifiers are vulnerable to gradient-based White-box attacks. And, even the model owner uses a defensively distilled model, the attacker can still have a chance to craft successful targeted attacks. We then proposed a modification to the standard DNN-based classifiers which helps to mask the gradients of the model and prevents the attacker to exploit them to craft both targeted and untargeted adversarial samples. We empirically verified the effectiveness of our approach on standard datasets which are heavily used by adversarial ML community. Finally, we demonstrated that our proposed model variants have inherent resistance to Deepfool attack thanks to the increased non-linearity at the output layer.

In this study, we focused on securing DNN based classifiers against evasion attacks. However, it is shown that previous defense approaches on adversarial robustness suffer from privacy preservation issues [31]. In the future, we plan to evaluate our proposed models against privacy related attack strategies, specifically membership inference attacks.

### VI. APPENDIX - A

Note: The common softmax layers are omitted for simplicity. For our proposed methods, we have applied Sigmoid and Tanh activation layers just after the final fully connected layers. The model architectures are available in the shared Github repository.

TABLE V: Model Architectures used in our experiments

Dataset	Layer Type	Layer Information
MNIST - Digit	Convolution (padding:1) + ReLU	3 × 3 × 32
	Convolution (padding:1) + ReLU	3 × 3 × 32
	Max Pooling	2 × 2
	Convolution (padding:1) + ReLU	3 × 3 × 64
	Convolution (padding:1) + ReLU	3 × 3 × 64
	Max Pooling	2 × 2
	Fully Connected + ReLU	3136 × 200
	Dropout	p : 0.5
	Fully Connected + ReLU	200 × 200
	Dropout	p : 0.5
Fully Connected	200 × 10	
CIFAR10	Convolution (Padding = 1) + ReLU	3 × 3 × 32
	Convolution (Padding = 1) + ReLU	3 × 3 × 64
	Max Pooling (Stride 2)	2 × 2
	Convolution (Padding = 1) + ReLU	3 × 3 × 128
	Convolution (Padding = 1) + ReLU	3 × 3 × 128
	Max Pooling (Stride 2)	2 × 2
	Convolution (Padding = 1) + ReLU	3 × 3 × 256
	Convolution (Padding = 1) + ReLU	3 × 3 × 256
	Dropout	p : 0.5
	Max Pooling (Stride 2)	2 × 2
	Fully Connected + ReLU	4096 × 1024
	Dropout	p : 0.5
	Fully Connected + ReLU	1024 × 256
	Dropout	p : 0.5
Fully Connected	256 × 10	

## VII. APPENDIX - B

The gradient derivation of the cross-entropy loss coupled with the softmax activation function is described in this part. This derivation was detailed for the first time in [36]. We will be using the derivation explained by Katzir et al. in [28] as it is.

Softmax Function Gradient Derivation:

Let  $K$  represents number of classes in training data,  $y = (y_0, y_1, \dots, y_{K-1})$  denotes the one-hot encoded label information,  $z_i$  denotes the  $i^{th}$  component of the logits layer output given some network input  $x$ . The probability estimate of the  $i^{th}$  class associated with the input by the softmax function is:

$$P_i = \frac{e^{z_i}}{\sum_{k=0}^{K-1} e^{z_k}} \quad (17)$$

$P_i$ 's derivative with respect to  $z_k$  can then be calculated as below:

$$\frac{\partial P_i}{\partial z_j} = \frac{\partial \left( \frac{e^{z_i}}{\sum_{k=0}^{K-1} e^{z_k}} \right)}{\partial z_j} \quad (18)$$

In the case of  $i = j$ , we get:

$$\begin{aligned} \frac{\partial P_i}{\partial z_j} &= \frac{\partial \left( \frac{e^{z_i}}{\sum_{k=0}^{K-1} e^{z_k}} \right)}{\partial z_j} = \frac{e^{z_i} \sum_{k=0}^{K-1} e^{z_k} - e^{z_i} e^{z_j}}{\left( \sum_{k=0}^{K-1} e^{z_k} \right)^2} \\ &= \frac{e^{z_i}}{\sum_{k=0}^{K-1} e^{z_k}} \cdot \frac{\left( \sum_{k=0}^{K-1} e^{z_k} \right) - e^{z_j}}{\sum_{k=0}^{K-1} e^{z_k}} = P_i(1 - P_j) \end{aligned} \quad (19)$$

Likewise, when  $i \neq j$ , we get:

$$\begin{aligned} \frac{\partial P_i}{\partial z_j} &= \frac{\partial \left( \frac{e^{z_i}}{\sum_{k=0}^{K-1} e^{z_k}} \right)}{\partial z_j} = \frac{0 - e^{z_i} e^{z_j}}{\left( \sum_{k=0}^{K-1} e^{z_k} \right)^2} \\ &= \frac{-e^{z_i}}{\sum_{k=0}^{K-1} e^{z_k}} \cdot \frac{e^{z_j}}{\sum_{k=0}^{K-1} e^{z_k}} = -P_i P_j \end{aligned} \quad (20)$$

When we combine the two previous results, we get:

$$\frac{\partial P_i}{\partial z_j} = \begin{cases} P_i(1 - P_j), & \text{if } i = j \\ -P_i P_j, & \text{if } i \neq j \end{cases} \quad (21)$$

The cross-entropy loss  $L$  for any input  $x$  is formulated as :

$$L = - \sum_{i=0}^{K-1} y_i \cdot \log(P_i) \quad (22)$$

Assuming 'log' as natural logarithm (ln) for simplicity, we may formulate the gradient of the cross-entropy loss with respect to the  $i^{th}$  logit as below:

$$\begin{aligned} \frac{\partial P_i}{\partial z_j} &= \frac{\partial \left( - \sum_{i=0}^{K-1} y_i \cdot \log(P_i) \right)}{\partial z_i} \\ &= - \sum_{k=0}^{K-1} y_k \cdot \frac{\partial \log(P_k)}{\partial z_i} = - \sum_{k=0}^{K-1} y_k \cdot \frac{\partial \log(P_k)}{\partial P_k} \cdot \frac{\partial P_k}{\partial z_i} \\ &= - \sum_{k=0}^{K-1} y_k \cdot \frac{1}{P_k} \cdot \frac{\partial P_k}{\partial z_i} \end{aligned} \quad (23)$$

Combining Cross-Entropy and Softmax Function Derivatives:

Knowing from Eq. 21 that the softmax derivative equation for the case when  $i = j$  differs from the other cases, we rearrange the loss derivative equation slightly to differentiate this case from the others:

$$\begin{aligned} \frac{\partial L}{\partial z_i} &= - \sum_{k=0}^{K-1} y_k \cdot \frac{1}{P_k} \cdot \frac{\partial P_k}{\partial z_i} \\ &= -y_i \cdot \frac{1}{P_i} \cdot \frac{\partial P_i}{\partial z_i} - \sum_{k \neq i} y_k \cdot \frac{1}{P_k} \cdot \frac{\partial P_k}{\partial z_i} \end{aligned} \quad (24)$$

We can now apply the derivative of softmax we derived before to obtain:

$$\begin{aligned} -y_i \cdot \frac{P_i(1 - P_i)}{P_i} - \sum_{k \neq i} \frac{y_k \cdot (-P_k P_i)}{P_k} &= -y_i + y_i P_i + \sum_{k \neq i} y_k P_i \\ &= P_i \left( y_i + \sum_{k \neq i} y_k \right) - y_i \end{aligned} \quad (25)$$

Luckily, because  $Y$  is the one-hot encoded actual label vector, we know that:

TABLE VI: CNN model parameters

	MNIST (Digit)				CIFAR-10			
	Normal Model	Distilled Model	Ours (Sigmoid)	Ours (Tanh)	Normal Model	Distilled Model	Ours (Sigmoid)	Ours (Tanh)
Opt.	Adam	Adam	Adam	Adam	Adam	Adam	Adam	Adam
LR	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001
Batch S.	128	128	128	128	128	128	128	128
Dropout	0.5	0.5	0.5	0.5	0.5	0.5	0.05	0.25
Epochs	20	20	20	20	50	50	50	50
Temp.	1	20	20	20	1	50	50	50

TABLE VII: Parameters that are used in BIM and Targeted BIM attacks:  $\alpha$  denotes the step size and  $i$  denotes # of steps for a perturbation budget  $\epsilon$ 

Dataset	Parameters	$l_p$ norm
MNIST Digit	$\epsilon = 0.1$ & $0.2$ , $\alpha = \epsilon \cdot 0.1$ , $i = 20$	$l_\infty$
CIFAR10	$\epsilon = 3/255$ & $6/255$ , $\alpha = \epsilon \cdot 0.1$ , $i = 20$	$l_\infty$

$$y_i + \sum_{k \neq i} y_k = \sum_{k=0}^{K-1} y_k = 1 \quad (26)$$

and therefore we finally end up with below expression as the derivative of the loss with respect to any logit:

$$\frac{\partial L}{\partial z_i} = P_i \left( y_i + \sum_{k \neq i} Y_k \right) - y_i = P_i - y_i \quad (27)$$

## REFERENCES

- [1] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," 2014.
- [2] O. F. Tuna, F. O. Catak, and M. T. Eskil, "Exploiting epistemic uncertainty of the deep learning models to generate adversarial samples," *Multimedia Tools and Applications*, vol. 81, no. 8, pp. 11479–11500, Mar 2022.
- [3] A. Ilyas, L. Engstrom, and A. Madry, "Prior convictions: Black-box adversarial attacks with bandits and priors," 2019.
- [4] O. F. Tuna, F. O. Catak, and M. T. Eskil, "Closeness and uncertainty aware adversarial examples detection in adversarial machine learning," *Computers and Electrical Engineering*, vol. 101, p. 107986, 2022.
- [5] D. Meng and H. Chen, "Magnet: A two-pronged defense against adversarial examples," 2017.
- [6] M. Sato, J. Suzuki, H. Shindo, and Y. Matsumoto, "Interpretable adversarial perturbation in input embedding space for text," 2018.
- [7] N. Carlini and D. Wagner, "Audio adversarial examples: Targeted attacks on speech-to-text," 2018.
- [8] S. G. Finlayson, H. W. Chung, I. S. Kohane, and A. L. Beam, "Adversarial attacks against medical deep learning systems," 2019.
- [9] C. Sitawarin, A. N. Bhagoji, A. Mosenia, M. Chiang, and P. Mittal, "Darts: Deceiving autonomous cars with toxic signs," 2018.
- [10] N. Morgulis, A. Kreines, S. Mendelowitz, and Y. Weisglass, "Fooling a real car with adversarial traffic signs," 2019.
- [11] Z. Zheng and P. Hong, "Robust detection of adversarial attacks by modeling the intrinsic properties of deep neural networks," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018.
- [12] X. Huang, D. Kroening, W. Ruan, J. Sharp, Y. Sun, E. Thamo, M. Wu, and X. Yi, "A survey of safety and trustworthiness of deep neural networks: Verification, testing, adversarial attack and defence, and interpretability," *Computer Science Review*, vol. 37, p. 100270, 2020.
- [13] F. O. Catak, S. Sivaslioglu, and K. Sahinbas, "A generative model based adversarial security of deep learning and linear classifier models," 2020.
- [14] A. Qayyum, M. Usama, J. Qadir, and A. Al-Fuqaha, "Securing connected autonomous vehicles: Challenges posed by adversarial machine learning and the way forward," *IEEE Communications Surveys Tutorials*, vol. 22, no. 2, pp. 998–1026, 2020.
- [15] K. Sadeghi, A. Banerjee, and S. K. S. Gupta, "A system-driven taxonomy of attacks and defenses in adversarial machine learning," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 4, no. 4, pp. 450–467, 2020.
- [16] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," 2015.
- [17] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," 2017.
- [18] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," 2016.
- [19] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," 2017.
- [20] M. Andriushchenko, F. Croce, N. Flammarion, and M. Hein, "Square attack: a query-efficient black-box adversarial attack via random search," 2020.
- [21] J. Chen, M. I. Jordan, and M. J. Wainwright, "Hopskipjumpattack: A query-efficient decision-based attack," in *2020 IEEE Symposium on Security and Privacy (SP)*, 2020, pp. 1277–1294.
- [22] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," 2015.
- [23] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," 2016.
- [24] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," 2019.
- [25] F. Liao, M. Liang, Y. Dong, T. Pang, X. Hu, and J. Zhu, "Defense against adversarial attacks using high-level representation guided denoiser," 2018.
- [26] S. Shen, G. Jin, K. Gao, and Y. Zhang, "Ape-gan: Adversarial perturbation elimination with gan," 2017.
- [27] A. Raghunathan, J. Steinhardt, and P. Liang, "Certified defenses against adversarial examples," 2020.
- [28] Z. Katzir and Y. Elovici, "Gradients cannot be tamed: Behind the impossible paradox of blocking targeted adversarial attacks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 1, pp. 128–138, 2021.
- [29] A. S. Ross and F. Doshi-Velez, "Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients," 2017.
- [30] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," 2017.
- [31] L. Song, R. Shokri, and P. Mittal, "Privacy risks of securing machine learning models against adversarial examples," *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, Nov 2019.
- [32] N. Carlini, A. Athalye, N. Papernot, W. Brendel, J. Rauber, D. Tsipras, I. Goodfellow, A. Madry, and A. Kurakin, "On evaluating adversarial robustness," 2019.
- [33] Y. LeCun and C. Cortes, "MNIST handwritten digit database," 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [34] A. Krizhevsky, V. Nair, and G. Hinton, "Cifar-10 (canadian institute for advanced research)." [Online]. Available: <http://www.cs.toronto.edu/~kriz/cifar.html>
- [35] J. Rauber, W. Brendel, and M. Bethge, "Foolbox: A python toolbox to benchmark the robustness of machine learning models," 2018.

- [36] D. Campbell, R. A. Dunne, and N. A. Campbell, "On the pairing of the softmax activation and cross-entropy penalty functions and the derivation of the softmax activation function."