

REGroup: Rank-aggregating Ensemble of Generative Classifiers for Robust Predictions

Lokender Tiwari^{1,2}Anish Madan¹Saket Anand¹Subhashis Banerjee^{3,4}

¹IIT-Delhi ²TCS Research ³IIT Delhi ⁴Department of Computer Science, Ashoka University
<https://lokender.github.io/REGroup.html>

Abstract

Deep Neural Networks (DNNs) are often criticized for being susceptible to adversarial attacks. Most successful defense strategies adopt adversarial training or random input transformations that typically require retraining or fine-tuning the model to achieve reasonable performance. In this work, our investigations of intermediate representations of a pre-trained DNN lead to an interesting discovery pointing to intrinsic robustness to adversarial attacks. We find that we can learn a generative classifier by statistically characterizing the neural response of an intermediate layer to clean training samples. The predictions of multiple such intermediate-layer based classifiers, when aggregated, show unexpected robustness to adversarial attacks. Specifically, we devise an ensemble of these generative classifiers that rank-aggregates their predictions via a Borda count-based consensus. Our proposed approach uses a subset of the clean training data and a pre-trained model, and yet is agnostic to network architectures or the adversarial attack generation method. We show extensive experiments to establish that our defense strategy achieves state-of-the-art performance on the ImageNet validation set.

1. Introduction

Deep Neural Networks (DNNs) have shown outstanding performance on many computer vision tasks such as image classification [27], speech recognition [22], and video classification [26]. Despite showing superhuman capabilities in the image classification task [21], the existence of *adversarial examples* [48] have raised questions on the reliability of neural network solutions for safety-critical applications.

Adversarial examples are carefully manipulated adaptations of an input, generated with the intent to fool a classifier into misclassifying them. One of the reasons for the attention that adversarial examples garnered is the ease with which they can be generated for a given model by simply maximizing the corresponding loss function. This is

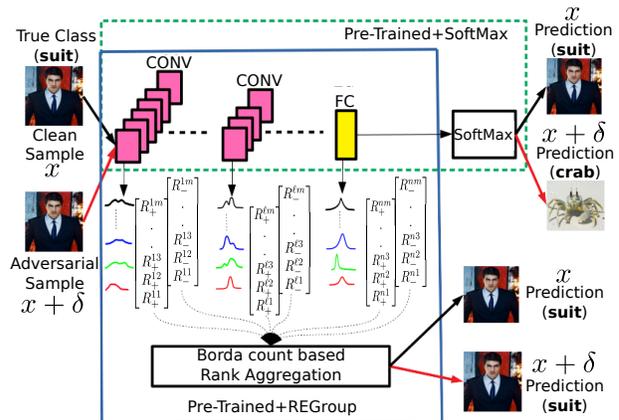


Figure 1. **Overview of REGroup.** Rank-aggregating Ensemble of Generative classifiers for robust predictions. REGroup uses a pre-trained network, and constructs layer-wise generative classifiers modeled by a mixture distribution of the positive and negative pre-activation neural responses at each layer. At test time, an input sample’s neural responses are tested with generative classifiers to obtain ranking preferences of classes at each layer. These preferences are aggregated using *Borda count* based preferential voting theory to make final prediction. *Note:* construction of layer-wise generative classifiers is a one time process.

achieved by simply using a gradient based approach that finds a small perturbation at the input which leads to a large change in the output [48]. This apparent instability in neural networks is most pronounced for deep networks that have an accumulation effect over the layers. This effect results in taking the small, additive, adversarial noise at the input and amplifying it to substantially noisy feature maps at intermediate layers that eventually influences the softmax probabilities enough to misclassify the perturbed input sample. This observation of amplification of input noise over the layers is not new, and has been pointed out in the past [48]. The recent work by [57] addresses this issue by introducing *feature denoising* blocks in a network and training them with adversarially generated examples.

The iterative nature of generating adversarial examples

makes their use in training to generate defenses computationally very intensive. For instance, the adversarially trained feature denoising model proposed by [57] takes 38 hours on 128 Nvidia V100 GPUs to train a baseline ResNet-101 with ImageNet. While we leverage this observation of noise amplification over the layers, our proposed approach *avoids any training or fine-tuning* of the model. Instead, we use a representative subset of training samples and their layer-wise pre-activation responses to construct mixture density based generative classifiers, which are then combined in an ensemble using ranking preferences.

Generative classifiers have achieved varying degrees of success as defense strategies against adversarial attacks. Recently, [18] studied the class-conditional generative classifiers and concluded that it is impossible to guarantee robustness of such models. More importantly, they highlight the challenges in training generative classifiers using maximum likelihood based objective and their limitations w.r.t. discriminative ability and identification of out-of-distribution samples. While we propose to use generative classifiers, we avoid using likelihood based measures for making classification decisions. Instead, we use rank-order preferences of these classifiers which are then combined using a *Borda count*-based voting scheme. Borda counts have been used in collective decision making and are known to be robust to various manipulative attacks [43].

In this paper, we present our defense against adversarial attacks on deep networks, referred to as *Rank-aggregating Ensemble of Generative classifiers for robust predictions* (REGroup). At inference time, our defense requires white-box access to a pre-trained model to collect the pre-activation responses at intermediate layers to make the final prediction. We use the training data to build our generative classifier models. Nonetheless, our strategy is simple, network-agnostic, does not require any training or fine-tuning of the network, and works well for a variety of adversarial attacks, even with varying degrees of hardness. Consistent with recent trends, we focus only on the ImageNet dataset to evaluate the robustness of our defense and report performance superior to defenses that rely on adversarial training [28] and random input transformation [41] based approaches. Finally, we present extensive analysis of our defense with two different architectures (ResNet and VGG) on different targeted and untargeted attacks. Our primary contributions are summarized below:

- We present REGroup, a retraining free, model-agnostic defense strategy that leverages an ensemble of generative classifiers over intermediate layers of the model.
- We model each layer-wise generative classifier as a simple mixture distribution of neural responses obtained from a subset of training samples. We discover

that *both positive and negative* pre-activation values contain information that can help correctly classify adversarially perturbed samples.

- We leverage the robustness inherent in Borda-count based consensus over the generative classifiers.
- We show extensive comparisons and analysis on the ImageNet dataset spanning a variety of adversarial attacks.

2. Related Work

Several defense techniques have been proposed to make neural networks robust to adversarial attacks. Broadly, we can categorize them into two approaches that: 1. Modify training procedure or modify input before testing; 2. Modify network or change hyper-parameters and optimization procedure.

2.1. Modify Training/Inputs During Testing

Some approaches of defenses in this category are mentioned below. *Adversarial training* [62, 34, 54, 45]. *Data compression* [4] suppresses the high-frequency components and presents an ensemble-based defense approach. *Data randomization* [53, 56] based approaches apply random transformations to the input to defend against adversarial examples by reducing their effectiveness.

PixelDefend [47] sets out to find the image with the highest probability within an ϵ -neighbourhood of the original image, thereby moving the image back towards distribution seen in training data. Defense-GAN [44] tries to model the distribution of unperturbed images and at inference, it generates an image close to what was provided but without adversarial perturbations. These two methods use techniques to generate a clean version of the input and pass to the classifier.

2.2. Modify Network/Network Add-ons

Defenses under this category address the *detection* of adversarial attacks or cater to both *detection and correction* of prediction. The aim of detection only defenses is to highlight if an example is adversarial and prevent it from further processing. These approaches include employing a detector sub-network [33], training the main classifier with an outlier class [19], using convolution filter statistics [29], or applying feature squeezing [58] to detect adversarial examples. However, all of these methods have shown to be ineffective against strong adversarial attacks [9][46]. Full defense approaches include applying defensive distillation [39][37] to use the knowledge from the output of the network to re-train the original model and improve the resilience of a network to small perturbations. Another approach is to augment the network with a sub-network called Perturbation Rectifying

Network (PRN) [1] to detect the perturbations; if the perturbation is detected, then PRN is used to classify the input image. However, later it was shown that the Carlini Wagner (C&W) attack successfully defeated the defensive distillation approach.

2.3. ImageNet Focused Defense Approaches

A few approaches have been evaluated on the ImageNet dataset, most of which are based on input transformations or image denoising. Nearly all these defenses designed for ImageNet have failed a thorough evaluation, with a regularly updated list maintained at [31]. The approaches in [40] and [30] claimed 81% and 75% accuracy respectively under adversarial attacks. But after a thorough evaluation [2] and accounting for obfuscated gradients [3], the accuracy for both was reduced to 0%. Similarly, [55] and [20] claimed 86% and 75% respectively, but these were also reduced to 0% [3]. A different approach proposed in [25] claimed an accuracy 27.9% but later it was also reduced to 0.1% [16]. For a comprehensive related work on attacks and defenses, we suggest reader to refer [11].

3. REGroup Methodology

Well-trained deep neural networks have a hierarchical structure, where the early layers transform inputs to feature spaces capturing local or more generic patterns, while later layers aggregate this local information to learn more semantically relevant representations. In REGroup, we use many of the higher layers and learn class-conditional generative classifiers, which are simple mixture-distributions estimated from the pre-activation neural responses at each layer from a subset of training samples. An ensemble of these layer-wise generative classifiers is used to make the final prediction by performing a Borda count-based rank-aggregation. Ranking preferences have been used extensively in robust fitting problems in computer vision [12, 23, 49], and we show its effectiveness in introducing robustness in DNNs against adversarial attacks.

Fig. 1 illustrates the overall working of REGroup. The approach has three main components: First, we use each layer as a generative classifier that produces a ranking preference over all classes. Second, each of these class-conditional generative classifiers are modeled using a mixture-distribution over the neural responses of the corresponding layer. Finally, the individual layer’s class ranking preferences are aggregated using Borda count-based scoring to make the final predictions. We introduce the notation below and discuss each of these steps in detail in the subsections that follow.

Notation. In this paper, we will always use ℓ , i and j for indexing the ℓ^{th} layer, i^{th} feature map and the j^{th} input sample respectively. The true and predicted class label will be denoted by y and \hat{y} respectively. A classifier can be represented

in a functional form as $\hat{y} = \mathcal{F}(x)$, it takes an input x and predicts its class label \hat{y} . We define $\phi^{\ell i}$ as the ℓ^{th} layer’s i^{th} pre-activation feature map, i.e., the neural responses *before* they pass through the activation function. For convolutional layers, this feature map $\phi^{\ell i}$ is a 2D array, while for a fully connected layer, it is a scalar value.

3.1. DNN Layers as Generative Classifiers

We use the highest k layers of a DNN as generative classifiers that use the pre-activation neural responses to produce a ranking preference¹ over all classes. The layer-wise generative classifiers are modeled as a class-conditional mixture distribution, which is estimated using only a *pre-trained* network and a small subset \mathcal{S} of the training data. Let \mathcal{S} contain only correctly classified training samples², which we can further divide into M subsets, one for each class i.e $\mathcal{S} = \{\cup_{y=1}^M \mathcal{S}_y\}$, where \mathcal{S}_y is the subset containing samples that have labels y .

3.1.1 Layerwise Neural Response Distributions

Our preliminary observations indicated that while the ReLU activations truncate the negative pre-activations during the forward pass, these values still contain semantically meaningful information. Our ablative studies in Fig. 5 confirm this observation and additionally, on occasion, we find that the negative pre-activations are complementary to the positive ones. Since the pre-activation features are real-valued, we compute the features $\phi_j^{\ell i}$ for the j^{th} sample x_j , and define its positive ($P_j^{\ell i}$) and negative ($N_j^{\ell i}$) response accumulators as $P_j^{\ell i} = \sum \max(0, \phi_j^{\ell i})$, $N_j^{\ell i} = \sum \max(0, -\phi_j^{\ell i})$.

For convolutional layers, these accumulators represent the overall strength of positive and negative pre-activation responses respectively, when aggregated over the spatial dimensions of the i^{th} feature map of the ℓ^{th} layer. On the other hand, for the linear layers, the accumulation becomes trivial with each neuron having a scalar response $\phi_j^{\ell i}$. We can now represent the ℓ^{th} layer by the positive and negative response accumulator vectors denoted by P_j^ℓ and N_j^ℓ respectively. We normalize these vectors and define the layer-wise probability mass function (PMF) for the positive and negative responses as $\mathbb{P}_j^\ell = \frac{P_j^\ell}{\|P_j^\ell\|_1}$ and $\mathbb{N}_j^\ell = \frac{N_j^\ell}{\|N_j^\ell\|_1}$ respectively.

Our interpretation of \mathbb{P}_j^ℓ and \mathbb{N}_j^ℓ as a PMF could be justified by drawing an analogy to the softmax output, which is also interpreted as a PMF. However, it is worth emphasizing that we chose the linear rescaling of the accumulator

¹A rank is assigned to each class based on a score. In the case of ImageNet dataset, the class with rank-1 is most preferred/likely class, while rank-1000 is the least preferred/likely class

²We took 50,000 out of ~ 1.2 millions training images from ImageNet dataset, 50 per class.

vectors rather than directly applying a softmax normalization. By separating out the positive and negative accumulators, we obtain two independent representations for each layer, which is beneficial to our rank-aggregating ensemble discussed in the following sections. A softmax normalization over a feature map comprising of positive and negative responses would have entirely suppressed the negative responses, discarding all its constituent semantic information. An additional benefit of the linear scaling is its simple computation. Algorithm 1 summarizes the computation of the layer-wise PMFs for a given training sample.

Algorithm 1: Layerwise PMF of neural responses. $H \times W$ represents the spatial dimensions of pre-activation features. For ℓ^{th} convolutional layer the dimensions of feature maps $H \times W = r^\ell \times s^\ell$, and for linear layers the dimensions of neuron output $H \times W = 1 \times 1$.

Input: x_j pre-activation features $\phi_j^{\ell i} \in \mathbb{R}^{H \times W}$
for $\ell \in [1..n]$ **do**
 $P_j^{\ell i} = \sum \max(0, \phi_j^{\ell i}), \quad \forall i$ (sum over H, W)
 $N_j^{\ell i} = \sum \max(0, -\phi_j^{\ell i}), \quad \forall i$ (sum over H, W)
end
 $P_j^\ell \leftarrow P_j^\ell + \delta, \quad N_j^\ell \leftarrow N_j^\ell + \delta$
 $\mathbb{P}_j^{\ell i} \leftarrow \frac{P_j^{\ell i}}{\sum_i \mathbb{P}_j^{\ell i}}, \quad \mathbb{N}_j^{\ell i} \leftarrow \frac{N_j^{\ell i}}{\sum_i \mathbb{N}_j^{\ell i}}$ (PMFs)

3.1.2 Layerwise Generative Classifiers

We model the layerwise generative classifiers for class y as a class-conditional mixture of distributions, with each mixture component as the PMFs \mathbb{P}_j^ℓ and \mathbb{N}_j^ℓ for a given training sample $x_j \in \mathcal{S}_y$. The generative classifiers corresponding to the positive and negative neural responses are then defined as the following mixture of PMFs

$$\mathbb{C}_y^{+\ell} = \sum_{j: x_j \in \mathcal{S}_y} \lambda_j \mathbb{P}_j^\ell, \quad \mathbb{C}_y^{-\ell} = \sum_{j: x_j \in \mathcal{S}_y} \lambda_j \mathbb{N}_j^\ell \quad (1)$$

where the weights λ_j are nonnegative and add up to one in the respective equations. Here, λ_j is proportional to the softmax probability of the sample x_j , and δ is the small constant used for numerical stability. We choose the weights to be proportional to the softmax probability value as predicted by the network given the input x_j . Using the subset of training samples \mathcal{S} , we construct the class-conditional mixture distributions, $\mathbb{C}_y^{+\ell}$ and $\mathbb{C}_y^{-\ell}$ at each layer ℓ only once. At inference time, we input a test sample x_j , from the test set \mathcal{T} , to the network and compute the PMFs \mathbb{P}_j^ℓ and \mathbb{N}_j^ℓ using Algorithm 1. As our test input is a PMF and

the generative classifier is also a mixture distribution, we simply use the KL-Divergence between the classifier model $\mathbb{C}^{+\ell}$ and the test sample \mathbb{P}_j^ℓ as a classification score as

$$P_{KL}(\ell, y) = \sum_i \mathbb{C}_y^{+\ell i} \log \left(\frac{\mathbb{C}_y^{+\ell i}}{\mathbb{P}_j^{\ell i}} \right), \forall y \in \{1, \dots, M\} \quad (2)$$

and similarly for the negative PMFs

$$N_{KL}(\ell, y) = \sum_i \mathbb{C}_y^{-\ell i} \log \left(\frac{\mathbb{C}_y^{-\ell i}}{\mathbb{N}_j^{\ell i}} \right), \forall y \in \{1, \dots, M\} \quad (3)$$

We use a simple classification rule and select the predicted class \hat{y} as the one with the smallest KL-Divergence with the test sample PMF. However, rather than identifying \hat{y} , at this stage we are only interested in rank-ordering the classes, which we simply achieve by sorting the KL-Divergences (Eqns. (2) and (3)) in ascending order. The resulting ranking preferences of classes for the ℓ^{th} layer are given below in Eqns. (4) and (5) respectively. Where, $R_+^{\ell y}$ is the rank (position of y^{th} class in the ascending order of KL-Divergences in P_{KL}) of y^{th} class in the ℓ^{th} layer preference list R_+^ℓ .

$$R_+^\ell = [R_+^{\ell 1}, R_+^{\ell 2}, \dots, R_+^{\ell y}, \dots, R_+^{\ell M}] \quad (4)$$

$$R_-^\ell = [R_-^{\ell 1}, R_-^{\ell 2}, \dots, R_-^{\ell y}, \dots, R_-^{\ell M}] \quad (5)$$

3.2. Robust Predictions with Rank Aggregation

Rank aggregation based preferential voting for making group decisions is widely used in selecting a winner in a democratic setup [43]. The basic premise of preferential voting is that n voters are allowed to rank m candidates in the order of their preferences. The rankings of all n voters are then aggregated to make a final prediction.

Borda count [6] is one of the approaches for preferential voting that relies on aggregating the rankings of all the voters to make a collective decision [43, 24]. The other popular voting strategies to find a winner out of m different choices include Plurality voting [52], and Condorcet winner [60]. In Plurality voting, the winner would be the one who gets the maximum fraction of votes, while Condorcet winner is the one who gets the majority votes.

3.2.1 Rank Aggregation Using Borda Count

Borda count is a generalization of the majority voting. In a two-candidates case it is equivalent to majority vote. The *Borda count* for a candidate is the sum of the number of candidates ranked below it by each voter. In our setting, while processing a test sample $x_j \in \mathcal{T}$, every layer acts as two independent voters based on \mathbb{P}^ℓ and \mathbb{N}^ℓ . The number of classes i.e M is the number of candidates. The Borda

count for the y^{th} class at the ℓ^{th} layer is denoted by $B^{\ell y} = B_+^{\ell y} + B_-^{\ell y}$, where $B_+^{\ell y}$ and $B_-^{\ell y}$ are the individual Borda count of both the voters and computed as shown in eq. (6).

$$B_+^{\ell y} = (M - R_+^{\ell y}), \quad B_-^{\ell y} = (M - R_-^{\ell y}) \quad (6)$$

3.2.2 Hyperparameter settings

We aggregate the Borda counts of highest k layers of the network, which is the only hyperparameter to set in REGroup. Let $B^{:ky}$ denote the aggregated Borda count of y^{th} class from the last k layers irrespective of the type (convolutional or fully connected). Here, n is the total number of layers. The final prediction would be the class with maximum aggregated Borda count.

$$\begin{aligned} B^{:ky} &= \sum_{\ell=n-k+1}^n B^{\ell y} \\ &= \sum_{\ell=n-k+1}^n B_+^{\ell y} + B_-^{\ell y}, \quad \forall y \in \{1..M\} \\ \hat{y} &= \operatorname{argmax}_y B^{:ky} \end{aligned} \quad (7)$$

To determine the value of k , we evaluate REGroup on 10,000 *correctly classified* samples from the ImageNet Validation set at each layer, using per layer Borda count i.e $\hat{y} = \operatorname{argmax}_y B^{\ell y}$. We select k to be the number of later layers at which we get at-least 75% accuracy. This can be viewed in the context of the confidence of individual layers on discriminating samples of different classes. We follow the above heuristic and found $k = 5$ for both the architectures ResNet-50 and VGG-19, which we use in all our experiments. An ablation study with all possible values of k is included in section 5.

4. Experiments

In this section, we evaluate robustness of REGroup against state-of-the-art attack methods. We follow the recommendations on defense evaluation in [8].

Attack methods. We consider attack methods in the following two categories: *gradient-based* and *gradient-free*.

Gradient-Based Attacks. Within this category, we consider two variants, *restricted* and *unrestricted* attacks. The restricted attacks generate adversarial examples by searching an adversarial perturbations within the bound of L_p norm, while unrestricted attacks generate adversarial example by manipulating image-based visual descriptors. Due to restriction on the perturbation the adversarial examples generated by restricted attacks are similar to the clean original image, while unrestricted attacks generate natural-looking adversarial examples, which are far from the clean original image in terms of L_p distance. We consider the following, *Restricted attacks*: PGD [32], DeepFool [35], C&W [10] and

Trust Region [59], and *Unrestricted attack*: cAdv [5] semantic manipulation attack. An example of cAdv is shown in Fig. 2.

Gradient-Free Attacks. The approaches in this category does not have access to the network weights. We consider following attacks: SPSA [51], Boundary [7] and Spatial [17]. Refer supplementary for the attack specific hyperparameters detail.



Figure 2. cAdv [5] adversarial examples

Network architectures. We consider ResNet-50³ and VGG-19⁴ architectures, pre-trained on ImageNet dataset.

Datasets. We present our evaluations, comparisons and analysis only on ImageNet [14] dataset. We use the subsets of full ImageNet validation set as described in Tab. 1. Note: V10K, V2K and V10C would be different for ResNet-50 and VGG-19, since an image classified correctly by ResNet-50 need not be classified correctly by the VGG-19.

Dataset	Description
V50K	Full ImageNet validation set with 50000 images.
V10K	A subset of 10000 correctly classified images from V50K set. 10 Per class.
V2K	A subset of 2000 correctly classified images from V50K set. 2 Per class.
V10C	A subset of correctly classified images of 10 sufficiently different classes.

Table 1. Dataset used for evaluation and analysis.

4.1. Performance on Gradient-Based Attacks

Comparison with adversarial-training/ fine-tuning. We evaluate REGroup on clean samples as well as adversarial examples generated using PGD ($\epsilon = 16$) from V50K dataset, and compare it with prior state-of-the-art works. The results are reported in Tab. 2, and we see that REGroup outperforms the state-of-the-art input transformation based defense BaRT [41], both in terms of the clean and adversarial samples (except in the case of Top-1 accuracy with $\hat{k} = 10$, which is the number of input transformations used in BaRT). We see that while our performance on clean samples decreases when compared to adversarial training (Inception v3), it improves significantly on adversarial examples with a high $\epsilon = 16$. While our method is not directly comparable with adversarially trained Inception v3 and ResNet-152, because the base models are different, a similar decrease in the accuracy over clean samples is reported in their paper. The trade-off between robustness and the standard accuracy has been studied in [15] and [50].

³<https://download.pytorch.org/models/resnet50-19c8e357.pth>

⁴<https://download.pytorch.org/models/vgg19-dcbb9e9d.pth>

An important observation to make with this experiment is, if we set aside the base models of ResNets and compare Top-1 accuracies on clean samples of full ImageNet validation set, our method (REGroup) without any *adv-training/fine-tuning* either outperforms or performs similar to the state-of-the-art *adv-training/fine-tuning* based methods [41, 57].

(Dataset used: V50K). Model	Clean Images		Attacked Images	
	Top-1	Top-5	Top-1	Top-5
ResNet-50	76	93	0.0	0.0
Inception v3	78	94	0.7	4.4
ResNet-152	79	94	-	-
Inception v3 w/Adv. Train	78	94	1.5	5.5
ResNet-152 w/Adv. Train	63	-	45	-
ResNet-152 w/Adv. Train w/ denoise	66	-	49	-
ResNet-50-BaRT, $\hat{k} = 5$	65	85	16	51
ResNet-50-BaRT, $\hat{k} = 10$	65	85	36	57
ResNet-50-REGroup	66	86	22	65

Table 2. Comparison with adversarially trained and fine-tuned classification models. Top-1 and Top-5 classification accuracy (%) of adversarial trained (Inception V3 [28] and ResNet-152 [57]) and fine-tuned (ResNet-50 BaRT [41]) classification models. Clean Images are the non-attacked original images. The results are divided into three blocks, the top block include original networks, middle block include defense approaches based on adversarial re-training/fine-tuning of original networks, bottom block is our defense *without re-training/fine-tuning*. Results of the competing methods are taken from their respective papers. ‘-’ indicate the results were not provided in the respective papers.

Performance w.r.t PGD Adversarial Strength. We evaluate REGroup w.r.t the maximum perturbation of the adversary. The results are reported in Fig. 3(a). REGroup outperforms both the adversarial training [28] and BaRT [41]. Both adversarial training and BaRT have shown protection against PGD adversarial attacks with a maximum perturbation strength $\epsilon = 16$ and $\epsilon = 32$ respectively, however we additionally show the results with $\epsilon = 40$ on full ImageNet validation set. We also note that with increasing perturbation strength, our defense’s accuracy is also strictly decreasing. This is in accordance with [8], where transitioning from a clean image to noise should yield a downward slope in accuracy, else there could be some form of gradient masking involved. While it may seem $\epsilon = 40$ is a large perturbation budget and it will destroy the object information in the image completely, but we would like to emphasize that it is not the case when using large size images. A comparison of PGD examples generated with $\epsilon = 40$ using CIFAR-10 (32×32) and ImageNet (224×224) images is shown in Fig. 3(b).

Performance on Un-Targeted Attacks. We evaluate REGroup on various untargeted attacks and report results in Tab. 3. The perturbation budgets (ϵ) and dataset used for the respective attacks are listed in the table. With the exception of the maximum perturbation allowed, we used

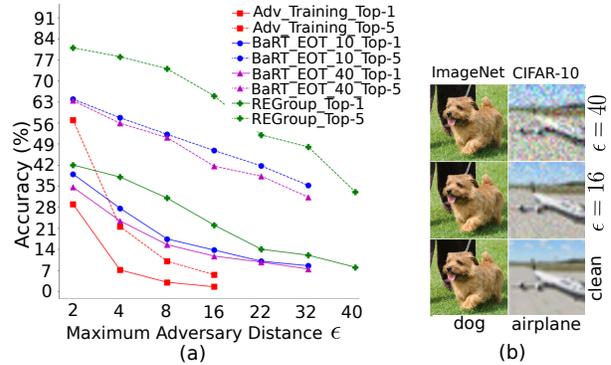


Figure 3. Top-1 and Top-5 accuracy(%) w.r.t PGD adversarial strength. Comparison with adversarial training based method [28] and fine-tuning using random input transformations based method (BaRT) [41] with Expectation Over Transformation (EOT) steps 10 and 40, against the PGD perturbation strength (ϵ). The results of the competing methods are taken from their respective papers. Dataset used: V50K.

	Data	UN / TA / HC	ϵ	ResNet-50			VGG-19		
				SMax		REGroup	SMax		REGroup
				#S	T1(%)	T1(%)	#S	T1(%)	T1(%)
Clean	V10K	-	-	10000	100	88	10000	100	76
Clean	V2K	-	-	2000	100	86	2000	100	72
Clean	V10C	-	-	417	100	84	392	100	79
PGD	V10K	UN	4 (L_∞)	9997	0	48	9887	0	46
DFool	V10K	UN	2 (L_2)	9789	0	61	9939	0	55
C&W	V10K	UN	4 (L_2)	10000	0	40	10000	0	38
TR	V10K	UN	2 (L_∞)	10000	0	41	9103	0	45
cAdv	V10C	UN	-	417	0	37	392	0	18
PGD	V2K	TA	(L_∞)	2000	0	47	2000	0	31
C&W	V2K	TA	(L_2)	2000	0	46	2000	0	38
PGD	V2K	UN+HC	(L_∞)	2000	0	21	2000	0	19
PGD	V2K	TA+HC	(L_∞)	2000	0	23	2000	0	17

Table 3. Performance on Gradient-Based Attacks. Comparison of Top-1 classification accuracy between SoftMax (SMax) and REGroup based final classification. UN and TA indicates, untargeted and targeted attacks respectively. The +HC indicates adversarial examples are generated with high-confidence ($> 90\%$) constraint, in this case ϵ can be any value that satisfies the HC criteria. For targeted attack we select a target class uniformly at random from the 1000 classes leaving out the true class. #S is the number of images for which the attacker is successfully able to generate adversarial examples using the respective attack models and the accuracies are reported with respect to the #S samples, hence the 0% accuracies with the SoftMax (SMax). Since #S is different for several attacks, therefore, the performance may not be directly comparable *across* different attacks. ‘-’ indicate the information is not-applicable. For data description refer Tab. 1.

default parameters given by FoolBox [42]. Due to space limitations, the attack specific hyper-parameters detail are included in the supplementary. We observe that the performance of our defense is quite similar for both the models employed. This is due to the attack-agnostic nature of our defense. We achieve 48% accuracy (ResNet-50) for PGD attack using our defense which is significant given that PGD is considered to be one of the strongest attacks

among the class of first order adversaries.

Performance on Unrestricted, Untargeted Semantic Manipulation Attacks. We consider V10C dataset for cAdv attack. We use the publicly released source code by the authors. Specifically we use cAdv₄ variant with the parameters suggested by the authors. The results are reported in Tab. 3.

Performance on Targeted Attacks. We consider V2K dataset for targeted attacks and report the performance on PGD and C&W targeted attacks in Tab. 3. Target class for such attacks is chosen uniformly at random from the 1000 ImageNet classes apart from the original(ground-truth) class.

Performance on PGD attack with High Confidence. We evaluate REGroup on PGD examples on which the network makes highly confident predictions using SoftMax. We generate un-targeted and targeted adversarial examples using PGD attack with a constraint that the network’s confidence of the prediction of adversarial examples is at-least 90%. For this experiment we do not put constraint on the adversarial perturbation i.e ϵ . Results are reported in Tab. 3.

4.2. Performance on Gradient-Free Attacks

Several studies [3], [38] have observed a phenomenon called *gradient masking*. This phenomenon occurs when a practitioner unintentionally or intentionally proposes a defense which does not have meaningful gradients, either by reducing them to small values (vanishing gradients), removing them completely (shattered gradients) or adding some noise to it (stochastic gradient). Gradient masking based defenses hinder the gradient computation and in turn inhibit gradient-based attacks, thus providing a false sense of security. Therefore, to establish the robustness of a defense against adversarial attacks in general, it is important to rule out that a defense relies on gradient masking.

To ensure that REGroup is not masking the gradients we follow the standard practice [36] [61] and evaluate on strong gradient-free SPSA [51] attack. In addition to SPSA, we also show results on two more gradient-free attacks, Boundary [7] and Spatial [17] attack. The results are reported in Tab. 4.

					ResNet-50			VGG-19		
	UN /		ϵ	#S	T1(%)	REGroup	#S	T1(%)	REGroup	
	Data	TA / HC								SMax
SPSA	V10K	UN	4 (L_∞)	4911	0	71	5789	0	58	
Boundary	V10K	UN	2 (L_2)	10000	0	50	10000	0	50	
Spatial	V10K	UN	2 (L_2)	2624	0	36	2634	0	30	

Table 4. **Performance on Gradient-Free Attacks.** Top-1 (%) classification accuracy comparison between SoftMax (SMax) and REGroup. Legends are same as in Tab. 3.

The consistent superior performance on both gradient-

based (both restricted and unrestricted) and gradient free attack shows REGroup is not masking the gradients and is attack method agnostic.

5. Analysis

5.1. Accuracy vs number of layers (k)

We report performance of REGroup on various attacks reported in Tab. 3 for all possible values of k . The accuracy of VGG-19 w.r.t. the various values of k is plotted in Fig. 4. We observe a similar accuracy vs k graph for ResNet-50 and note that a reasonable choice of k made based on this graph does not significantly impact REGroup’s performance. Refer Fig. 4, the ‘Agg’ stands for using aggregated Borda count $B^{:ky}$. PGD(V10K,UN), DFool, C&W(V10K,UN) and Trust Region are the same experiments as reported in Tab. 3, but with all possible values of k . ‘Per_Layer_V10K’ stands for evaluation using per layer Borda count i.e $\hat{y} = \text{argmax}_y B^{\ell y}$ on a separate 10,000 correctly classified subset of validation set. In all our experiments we choose the k -highest layers where ‘Per_Layer_V10K’ has at-least 75% accuracy. A reasonable change in this accuracy criteria of 75% would not affect the results on adversarial attacks significantly. However, a substantial change (to say 50%) deteriorates the performance on clean sample significantly.

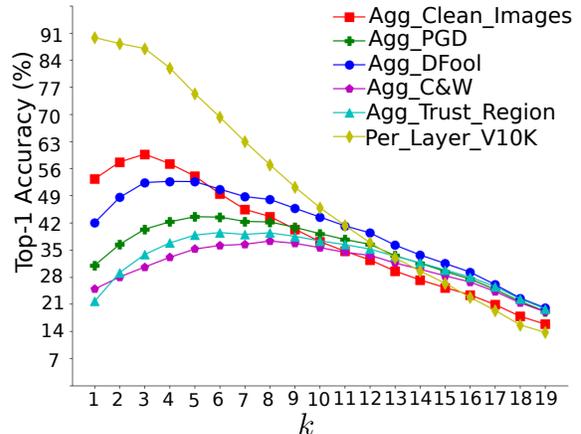


Figure 4. Accuracy vs no. of layers (k)

The phenomenon of decrease in accuracy of clean samples vs robustness has been studied in [15] and [50].

5.2. Effect of positive and negative pre-activation responses

We report the impact of using positive, negative and a combination of both pre-activation responses on the performance of REGroup in Fig. 5. We consider three variants of Borda count rank aggregation from later k layers. Pos: $B^{:ky} = \sum_{\ell=n-k+1}^n B_+^{\ell y}$, Neg: $B^{:ky} = \sum_{\ell=n-k+1}^n B_-^{\ell y}$, and Pos+Neg: $B^{:ky} = \sum_{\ell=n-k+1}^n B_+^{\ell y} + B_-^{\ell y}$. We report

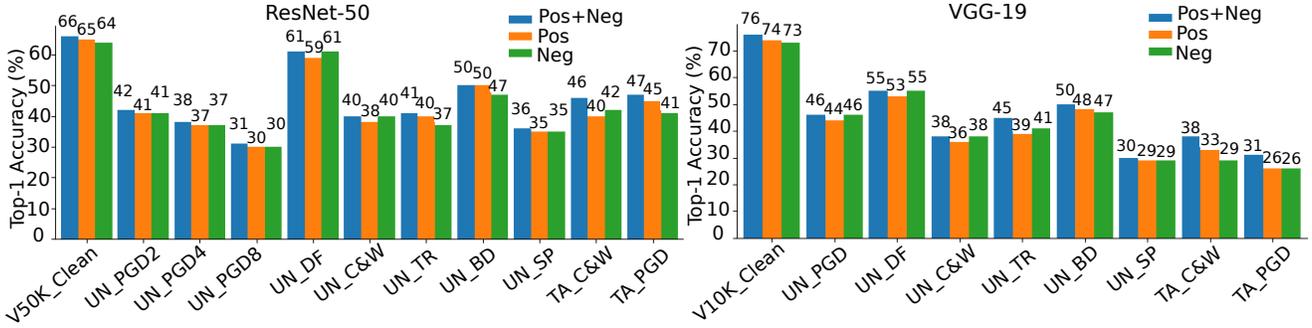


Figure 5. Effect of Considering Positive and Negative Pre-Activation Responses

the Top-1 accuracy (%) of the attacks experiment as set up in Tab. 3 (DF: DFool, C&W, TR: Trust Region), in Tab. 4 (BD: Boundary, SP: Spatial), and in Fig. 3 (PGD2, PGD4 and PGD8, with $\epsilon = 2, 4$ and 8 respectively). From the bar chart it is evident that in some experiments, Pos performs better than Neg (e.g UN_TR), while in others Neg is better than Pos only (e.g UN_DF). It is also evident that Pos+Neg occasionally improve the overall performance, and the improvement seems significant in the targeted C&W attacks for both the ResNet-50 and VGG-19. We leave it to the design choice of the application, if inference time is an important parameter, then one may choose either Pos or Neg to reduce the inference time to approximately half of what is reported in Tab. 6.

5.3. Results on CIFAR-10

While we mainly show results on large-scale dataset (ImageNet), we believe scaling down the datasets to one like CIFAR10 will not have a substantial impact on REGGroup’s performance. We evaluate REGGroup on CIFAR10 dataset using VGG-19 based classifier. We construct generative classifiers using CIFAR-10 dataset following the same protocol as for the ImageNet case described in the Sec. 3.1. We apply PGD attack with $\epsilon = 4$ and generate adversarial examples. The results are included in the Tab. 5.

	VGG-19		
	#S	SMax	REGGroup
		T1(%)	T1(%)
Clean	10000	92	88
PGD Untargeted	9243	0	57

Table 5. **Performance on CIFAR10.** Comparison of Top-1 classification accuracy between SoftMax (SMax) and REGGroup based final classification. #S is the number of images for which the attacker is successfully able to generate adversarial examples using PGD attack and the accuracies are reported with respect to the #S samples, hence the 0% accuracies with the SoftMax (SMax).

5.4. Inference time using REGGroup

Since we suggest to use REGGroup as a test time replacement of SoftMax, we compare the inference time on both

CPU and GPU in Tab. 6. We use a machine with an i7-8700 CPU and GTX 1080 GPU.

	ResNet-50				VGG-19			
	SMax		REGGroup		SMax		REGGroup	
	GPU	CPU	GPU	CPU	GPU	CPU	GPU	CPU
Time(s)	0.02	0.06	0.13	0.35	0.03	0.12	0.16	0.64

Table 6. **Inference time comparison.** REGGroup vs SoftMax

6. Conclusion

In this work, we have presented a simple, scalable, and practical defense strategy that is model agnostic and does not require any re-training or fine-tuning. We suggest to use REGGroup at test time to make a pre-trained network robust to adversarial perturbations. There are three aspects of REGGroup that justify its success. First, instead of using a maximum likelihood based prediction, REGGroup adopts a ranking preference based approach. Second, aggregation of preferences from multiple layers lead to group decision making, unlike SoftMax that relies on the output of the last layer only. Using both positive and negative layerwise responses help contribute to the robustness of REGGroup. Third, there exists inherent robustness of Borda count based rank aggregation in the presence of noisy individual voters [43], [24]. Hence, where SoftMax fails to predict the correct class of an adversarial example, REGGroup takes ranked preferences from multiple layers and builds a consensus using Borda count to make robust predictions. Our promising empirical results indicate that deeper theoretical analysis of REGGroup would be an interesting direction to pursue. One direction of analysis could be inspired from the recently proposed perspective of neurons as cooperating classifiers [13].

Acknowledgement

This work was partially supported by the Infosys Center for Artificial Intelligence, IIIT-Delhi.

References

- [1] Naveed Akhtar, Jian Liu, and Ajmal Mian. Defense against universal adversarial perturbations. In *CVPR*, pages 3389–3398, 2018.
- [2] Anish Athalye and Nicholas Carlini. On the robustness of the cvpr 2018 white-box adversarial example defenses. *arXiv preprint arXiv:1804.03286*, 2018.
- [3] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *ICML*, pages 274–283, 2018.
- [4] Arjun Nitin Bhagoji, Daniel Cullina, Chawin Sitawarin, and Prateek Mittal. Enhancing robustness of machine learning systems via data transformations. In *2018 52nd Annual Conference on Information Sciences and Systems (CISS)*, pages 1–5. IEEE, 2018.
- [5] Anand Bhattad, Min Jin Chong, Kaizhao Liang, Bo Li, and DA Forsyth. Unrestricted adversarial examples via semantic manipulation. In *ICLR*, 2020.
- [6] Duncan Black et al. The theory of committees and elections. 1958.
- [7] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *arXiv preprint arXiv:1712.04248*, 2017.
- [8] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, and Aleksander Madry. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.
- [9] Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pages 3–14. ACM, 2017.
- [10] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.
- [11] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. Adversarial attacks and defences: A survey. *arXiv preprint arXiv:1810.00069*, 2018.
- [12] Tat-Jun Chin, Jin Yu, and David Suter. Accelerated hypothesis generation for multistructure data via preference analysis. *IEEE TPAMI*, 34(4):625–638, 2011.
- [13] Marelle Davel, Marthinus Theunissen, Arnold Pretorius, and Etienne Barnard. Dnns as layers of cooperating classifiers. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):3725–3732, Apr 2020.
- [14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [15] Elvis Dohmatob. Limitations of adversarial robustness: strong no free lunch theorem. *arXiv preprint arXiv:1810.04065*, 2018.
- [16] Logan Engstrom, Andrew Ilyas, and Anish Athalye. Evaluating and understanding the robustness of adversarial logit pairing. *arXiv preprint arXiv:1807.10272*, 2018.
- [17] Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring the landscape of spatial robustness. In *ICML*, pages 1802–1811, 2019.
- [18] Ethan Fetaya, Joern-Henrik Jacobsen, Will Grathwohl, and Richard Zemel. Understanding the limitations of conditional generative models. In *ICLR*, 2020.
- [19] Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick McDaniel. On the (statistical) detection of adversarial examples. *arXiv preprint arXiv:1702.06280*, 2017.
- [20] Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens Van Der Maaten. Countering adversarial images using input transformations. *arXiv preprint arXiv:1711.00117*, 2017.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, pages 1026–1034, 2015.
- [22] Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdelrahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Brian Kingsbury, et al. Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal processing magazine*, 29, 2012.
- [23] Tat jun Chin, Hanzi Wang, and David Suter. The ordered residual kernel for robust motion subspace clustering. In Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, editors, *NIPS*, pages 333–341. Curran Associates, Inc., 2009.
- [24] Anson Kahng, Min Kyung Lee, Ritesh Noothigattu, Ariel Procaccia, and Christos-Alexandros Psomas. Statistical foundations of virtual democracy. In *ICML*, pages 3173–3182, 2019.
- [25] Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. *arXiv preprint arXiv:1803.06373*, 2018.
- [26] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, pages 1725–1732, 2014.
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012.
- [28] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*, 2016.
- [29] Xin Li and Fuxin Li. Adversarial examples detection in deep networks with convolutional filter statistics. In *ICCV*, pages 5764–5772, 2017.
- [30] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu. Defense against adversarial attacks using high-level representation guided denoiser. In *CVPR*, pages 1778–1787, 2018.
- [31] Aleksander Madry, Athalye Anish, Tsipras Dimitris, and Engstrom Logan. <https://www.robust-ml.org/>. <https://www.robust-ml.org/>, 2020. Accessed: 2020-01-05.
- [32] 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.

- [33] Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. *arXiv preprint arXiv:1702.04267*, 2017.
- [34] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *CVPR*, pages 1765–1773, 2017.
- [35] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *CVPR*, pages 2574–2582, 2016.
- [36] Tianyu Pang, Kun Xu, Yinpeng Dong, Chao Du, Ning Chen, and Jun Zhu. Rethinking softmax cross-entropy loss for adversarial robustness. In *ICLR*, 2020.
- [37] Nicolas Papernot and Patrick McDaniel. Extending defensive distillation. *arXiv preprint arXiv:1705.05264*, 2017.
- [38] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia conference on computer and communications security*, pages 506–519, 2017.
- [39] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *2016 IEEE Symposium on Security and Privacy (SP)*, pages 582–597. IEEE, 2016.
- [40] Aaditya Prakash, Nick Moran, Solomon Garber, Antonella DiLillo, and James Storer. Deflecting adversarial attacks with pixel deflection. In *CVPR*, pages 8571–8580, 2018.
- [41] Edward Raff, Jared Sylvester, Steven Forsyth, and Mark McLean. Barrage of random transforms for adversarially robust defense. In *CVPR*, pages 6528–6537, 2019.
- [42] 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.
- [43] Jörg Rothe. Borda count in collective decision making: A summary of recent results. In *AAAI*, volume 33, pages 9830–9836, 2019.
- [44] Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-gan: Protecting classifiers against adversarial attacks using generative models. *arXiv preprint arXiv:1805.06605*, 2018.
- [45] Ali Shafahi, Mahyar Najibi, Zheng Xu, John P Dickerson, Larry S Davis, and Tom Goldstein. Universal adversarial training. In *AAAI*, pages 5636–5643, 2020.
- [46] Yash Sharma and Pin-Yu Chen. Bypassing feature squeezing by increasing adversary strength. *arXiv preprint arXiv:1803.09868*, 2018.
- [47] Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *ICLR*, 2018.
- [48] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [49] Lokender Tiwari and Saket Anand. Dgsac: Density guided sampling and consensus. In *IEEE WACV*, pages 974–982. IEEE, 2018.
- [50] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- [51] Jonathan Uesato, Brendan O’Donoghue, Pushmeet Kohli, and Aaron Oord. Adversarial risk and the dangers of evaluating against weak attacks. In *ICML*, pages 5025–5034, 2018.
- [52] Jill Van Newenhizen. The borda method is most likely to respect the condorcet principle. *Economic Theory*, 2(1):69–83, 1992.
- [53] Qinglong Wang, Wenbo Guo, Kaixuan Zhang, II Ororbia, G Alexander, Xinyu Xing, Xue Liu, and C Lee Giles. Learning adversary-resistant deep neural networks. *arXiv preprint arXiv:1612.01401*, 2016.
- [54] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan L Yuille, and Quoc V Le. Adversarial examples improve image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 819–828, 2020.
- [55] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. *arXiv preprint arXiv:1711.01991*, 2017.
- [56] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. Adversarial examples for semantic segmentation and object detection. In *ICCV*, pages 1369–1378, 2017.
- [57] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *CVPR*, pages 501–509, 2019.
- [58] Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. *arXiv preprint arXiv:1704.01155*, 2017.
- [59] Zhewei Yao, Amir Gholami, Peng Xu, Kurt Keutzer, and Michael W Mahoney. Trust region based adversarial attack on neural networks. In *CVPR*, pages 11350–11359, 2019.
- [60] H Peyton Young. Condorcet’s theory of voting. *American Political science review*, 82(4):1231–1244, 1988.
- [61] Zhanyuan Zhang, Benson Yuan, Michael McCoyd, and David Wagner. Clipped bagnet: Defending against sticker attacks with clipped bag-of-features. In *3rd Deep Learning and Security Workshop (DLS)*, 2020.
- [62] Stephan Zheng, Yang Song, Thomas Leung, and Ian Goodfellow. Improving the robustness of deep neural networks via stability training. In *CVPR*, pages 4480–4488, 2016.