

σ -ZERO: GRADIENT-BASED OPTIMIZATION OF ℓ_0 -NORM ADVERSARIAL EXAMPLES

Anonymous authors

Paper under double-blind review

ABSTRACT

Evaluating the adversarial robustness of deep networks to gradient-based attacks is challenging. While most attacks focus ℓ_2 -norm and ℓ_∞ -norm constraints to craft input perturbations, only a few have investigated sparse ℓ_1 -norm and ℓ_0 -norm attacks. In particular, ℓ_0 -norm attacks remain the least studied due to the inherent complexity of optimizing over a non-convex and non-differentiable constraint. However, evaluating the robustness of these attacks might unveil weaknesses otherwise left untested with conventional ℓ_2 and ℓ_∞ attacks. In this work, we propose a novel ℓ_0 -norm attack, called σ -zero, which leverages an ad-hoc differentiable approximation of the ℓ_0 norm to facilitate gradient-based optimization. Extensive evaluations on MNIST, CIFAR10, and ImageNet datasets, involving robust and non-robust models, show that σ -zero can find minimum ℓ_0 -norm adversarial examples without requiring any time-consuming hyperparameter tuning, and that it outperforms all competing attacks in terms of success rate and scalability.

1 INTRODUCTION

Early research has unveiled that Deep Neural Networks (DNNs) are fooled by adversarial examples, i.e., slightly-perturbed inputs optimized to cause misclassifications (Biggio et al., 2013; Szegedy et al., 2014a; Goodfellow et al., 2015). In turn, this has demanded the need for more careful reliability assessments of such models. Most of the gradient-based attacks proposed to evaluate adversarial robustness of DNNs optimize adversarial examples under different ℓ_p -norm constraints. In particular, while convex ℓ_1 , ℓ_2 , and ℓ_∞ norms have been widely studied (Chen et al., 2018; Croce & Hein, 2021a), only a few ℓ_0 -norm attacks have been considered so far. The main reason is that ad-hoc heuristics need to be adopted to compute efficient projections on the ℓ_0 norm, overcoming issues related to its non-convexity and non-differentiability. Although this task is challenging and computationally expensive, attacks based on the ℓ_0 norm have the potential to reveal uncovered issues in DNNs that may not be evident in other norm-based attacks (Carlini & Wagner, 2017; Croce & Hein, 2021a). For instance, these attacks, known for perturbing a minimal fraction of input features, can be used to determine the most sensitive characteristics that influence the model’s decision-making process. Furthermore, they offer a different and relevant threat model to benchmark existing defenses. Developing efficient algorithms for generating ℓ_0 adversarial examples is thus a crucial area of research that requires further exploration to improve current adversarial robustness evaluations.

Unfortunately, current implementations of ℓ_0 attacks exhibit a largely suboptimal tradeoff between their success rate and efficiency, i.e., they are either accurate but slow, or fast but inaccurate. In particular, the accurate ones resort to the use of complex projections to find smaller input perturbations but suffer from time or memory limitations, hindering their scalability to larger networks or high-dimensional data (Brendel et al., 2019; Césaire et al., 2021). Other attacks execute faster, but their output solution is typically inaccurate and largely suboptimal as they rely on heuristic approaches and imprecise approximations to bypass the difficulties of optimizing the ℓ_0 norm, leading to overestimating adversarial robustness (Matyasko & Chau, 2021; Pintor et al., 2021). However, all existing strategies are often slow to converge because they require a large number of queries (i.e., forward and backward passes), or they output suboptimal solutions. It thus remains an open challenge to develop a scalable and compelling method for assessing the robustness of DNNs against sparse perturbations with minimum ℓ_0 norm.

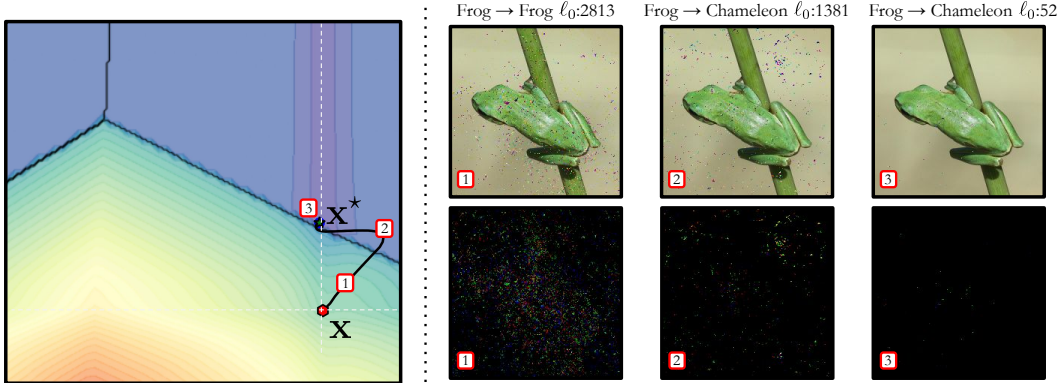


Figure 1: The leftmost plot shows an instance of σ -zero’s execution on a two-dimensional problem. The red dot and the green star respectively represent the initial point \mathbf{x} and the corresponding adversarial example \mathbf{x}^* . Our gradient-based attack seeks to find this adversarial example while minimizing the number of perturbed features (i.e., the ℓ_0 norm of the perturbation). Gray lines surrounding \mathbf{x} demarcate regions where the ℓ_0 norm is minimized. The rightmost plot shows the adversarial images (top row) and the corresponding perturbations (bottom row) found by σ -zero during the three steps highlighted in the leftmost plot, alongside their prediction and ℓ_0 norm.

To tackle these issues, in this work we propose a novel attack technique, namely σ -zero, which iteratively promotes the sparsity of the adversarial perturbations, minimizing their ℓ_0 norm (see Fig. 1 and Sect. 2). The underlying idea is to leverage a differentiable approximation of the actual ℓ_0 norm, which is better suited to gradient-based optimizers. Specifically, we employ the approximation initially introduced by Osborne et al. (2000b), and more recently adopted by Cinà et al. (2022) for staging energy-latency poisoning attacks. This method offers an unbiased, differentiable estimation of the true ℓ_0 norm, allowing us to optimize it via gradient descent.

Our experiments (Sect. 3) provide compelling evidence of the remarkable performance of our attack. We evaluate σ -zero on several benchmark datasets, including MNIST, CIFAR10, and ImageNet, considering baseline and robust models from Robustbench (Croce et al., 2021). We compare its performance with state-of-the-art attacks, showing that σ -zero achieves better results in terms of attack success rate and perturbation size, while being significantly faster and without requiring any sophisticated and time-consuming hyperparameter tuning. Overall, our approach encompasses two fundamental characteristics for a proficient adversarial attack, i.e., effectiveness and scalability, making it a catalyst for significant advancements in developing novel models with improved robustness, as well as better robustness evaluation tools.

2 σ -ZERO: MINIMUM ℓ_0 -NORM ADVERSARIAL EXAMPLES

We present here σ -zero, our gradient-based approach to finding minimum ℓ_0 -norm adversarial examples. We start by describing the considered threat model and then give a formal overview of the proposed attack and its algorithmic implementation.

Threat Model. We assume that the attacker has complete access to the target model, including its architecture and trained parameters, and exploits its gradient for staging white-box untargeted attacks. This setting is useful for worst-case evaluation of the adversarial robustness of DNN models, providing empirical upper bounds on the performance degradation that may incur when they are attacked, and it is the usual setting adopted also in previous work related to gradient-based adversarial robustness evaluations (Carlini & Wagner, 2017; Croce et al., 2021; Pintor et al., 2021).

Problem Formulation. In this work, we seek untargeted minimum ℓ_0 -norm adversarial perturbations that steer the model’s decision towards misclassification. To this end, let $\mathbf{x} \in \mathcal{X} = [0, 1]^d$ be a d -dimensional input sample, $y \in \mathcal{Y} = \{1, \dots, l\}$ its associated true label, and $f : \mathcal{X} \times \Theta \mapsto \mathcal{Y}$ the target model, parameterized by $\theta \in \Theta$. While f outputs the predicted label, we will also use f_k to denote the continuous-valued output (logit) for class $k \in \mathcal{Y}$. The goal of our attack is to find the

minimum ℓ_0 -norm adversarial noise δ^* such that the corresponding adversarial example $\mathbf{x}^* = \mathbf{x} + \delta^*$ is misclassified by f . This is formalized as the following optimization problem:

$$\delta^* \in \arg \min_{\delta} \quad \|\delta\|_0, \quad (1)$$

$$\text{s.t.} \quad f(\mathbf{x} + \delta, \theta) \neq y, \quad (2)$$

$$\mathbf{x} + \delta \in [0, 1]^d, \quad (3)$$

where $\|\cdot\|_0$ denotes the ℓ_0 norm, which counts the number of non-zero dimensions. The hard-constraint in Equation 2 ensures that the perturbation δ induces the target model f to misclassify the perturbed sample \mathbf{x}^* . Finally, Equation 3 represents a box constraint, ensuring that the adversarial example \mathbf{x}^* lies in $[0, 1]^d$. Note that when the source point \mathbf{x} is already misclassified by f , the trivial solution to the above minimization problem is $\delta^* = \mathbf{0}$.

Contrary to the $\ell_1, \ell_2, \ell_\infty$ norms, when considering the ℓ_0 norm the problem becomes intractable with standard methods. The ℓ_0 norm is indeed non-differentiable, thus unsuitable for gradient-based optimization. To address this issue, we exploit the ℓ_0 -norm approximation function proposed by Osborne et al. (2006b), and defined as:

$$\hat{\ell}_0(\mathbf{x}) = \sum_{i=1}^d \frac{x_i^2}{x_i^2 + \sigma}, \quad \sigma > 0, \quad \hat{\ell}_0(\mathbf{x}) \in [0, d], \quad (4)$$

with σ being a hyperparameter controlling its approximation quality. When σ tends to zero, the approximation becomes more accurate. However, an increasingly accurate approximation could lead to the same optimization limits of the ℓ_0 norm.

Finally, similarly to previous work (Carlini & Wagner, 2017; Rony et al., 2021a; Szegedy et al., 2014b), we transform the hard-constraint in Equation 2 in a soft-constraint. The resulting optimization problem therefore becomes:

$$\delta^* \in \arg \min_{\delta} \quad \mathcal{L}(\mathbf{x} + \delta, y, \theta) + \frac{1}{d} \hat{\ell}_0(\delta) \quad (5)$$

$$\text{s.t.} \quad \mathbf{x} + \delta \in [0, 1]^d, \quad (6)$$

where we substituted the $\|\delta\|_0$ with the approximation $\hat{\ell}_0(\delta)$ and normalize it with respect to the number of features d to ensure that its value is within the interval $[0, 1]$. The loss \mathcal{L} is defined as:

$$\mathcal{L}(\mathbf{x}, y, \theta) = \max \left(f_y(\mathbf{x}, \theta) - \max_{k \neq y} f_k(\mathbf{x}, \theta), 0 \right) + \mathbb{I}(f(\mathbf{x}, \theta) = y). \quad (7)$$

The first term in \mathcal{L} represents the logit difference, which is positive when the sample is correctly assigned to the true class y , and clipped to zero when it is misclassified (Carlini & Wagner, 2017). The second term merely adds 1 to the loss if the sample is correctly classified.¹ This ensures that the loss term \mathcal{L} is 0 only when an adversarial example is found, and higher than 1 otherwise. This in turn implies that the loss term \mathcal{L} is always higher than the ℓ_0 -norm term in Equation 5 (as the latter is bounded in $[0, 1]$), when no adversarial example is found. Accordingly, it is not difficult to see that the feasible solutions of this problem only correspond to minimum-norm adversarial examples. It is also worth remarking that, conversely to the objective function proposed by Carlini & Wagner (2017), our objective does not require tuning the tradeoff between minimizing the loss and reducing the perturbation size to find minimum-norm adversarial examples, thereby avoiding a computationally-expensive line search for each input sample. In fact, the proposed objective function inherently induces an *alternate* optimization process between the loss term and the ℓ_0 -norm penalty, as shown in the Appendix (see Figure 4). In particular, when the sample is not adversarial, the attack algorithm mostly aims to decrease the loss term to find an adversarial example, while increasing the perturbation size. Conversely, when an adversarial example is found, the loss term is cropped to zero, and the perturbation size is gradually reduced.

Solution Algorithm. Given that the approximation function $\hat{\ell}_0$ in Equation 4 is differentiable, we derive a custom gradient-based algorithm for solving Equation 5 and Equation 6. Our attack, detailed

¹While a sigmoid approximation may be adopted to overcome the non-differentiability of the \mathbb{I} term at the decision boundary, we simply set its gradient to zero *everywhere*, without any impact on the experimental results.

Algorithm 1: σ -zero Attack Pseudocode.

Input: $\mathbf{x} \in [0, 1]^d$, input sample; y , true class label; θ , target model; N , number of iterations;
 σ , $\hat{\ell}_0$ -approximation parameter; η_0 , initial step size; τ_0 , initial sparsity threshold.

Output: \mathbf{x}^* , minimum ℓ_0 norm adversarial example.

```

1  $\delta \leftarrow \mathbf{0}$ ;  $\delta^* \leftarrow \delta$ ;  $\tau \leftarrow \tau_0$ ;  $\eta \leftarrow \eta_0$  ▷ initialization.
2 for  $i$  in  $1, \dots, N$  do
3    $\nabla \mathbf{g} \leftarrow \nabla_{\delta} [\mathcal{L}(\mathbf{x} + \delta, y, \theta) + \frac{1}{d} \hat{\ell}_0(\delta, \sigma)]$  ▷ gradient computation.
4    $\nabla \mathbf{g} \leftarrow \nabla \mathbf{g} / \|\nabla \mathbf{g}\|_{\infty}$  ▷ gradient normalization.
5    $\delta \leftarrow \text{clip}(\mathbf{x} - [\delta - \eta \cdot \nabla \mathbf{g}]) - \mathbf{x}$  ▷  $\delta$  update.
6    $\delta \leftarrow \Pi_{\tau}(\delta)$  ▷ zeroing  $\delta$  components below  $\tau$ .
7    $\eta = \text{cosine\_annealing}(\eta_0, i)$  ▷  $\eta$  update.
8   if  $\mathcal{L}(\mathbf{x} + \delta, y, \theta) \leq 0$   $\tau + = 0.01 \cdot \eta$  else  $\tau - = 0.01 \cdot \eta$  ▷  $\tau$  update.
9 end
10 if  $\mathcal{L}(\mathbf{x} + \delta, y, \theta) \leq 0 \wedge \|\delta\|_0 < \|\delta^*\|_0$   $\delta^* \leftarrow \delta$  ▷  $\delta^*$  update.
11 return  $\mathbf{x}^* \leftarrow \mathbf{x} + \delta^*$ 

```

in Algorithm 1 is fast, not memory-demanding, and easy to implement. It starts by initializing the adversarial perturbation $\delta = \mathbf{0}$ (line 1). Subsequently, it computes the gradient of the objective function in Equation 5 with respect to δ (line 3), and normalizes it to speed up convergence (Rony et al., 2018; Pintor et al., 2021). We then update δ to minimize the objective via gradient descent, while also accounting for the box constraints in Equation 6 through the usage of the clip operator (line 5). We enforce sparsity in δ by clipping to 0 all the components lower than the current sparsity threshold τ (Line 6). This step is necessary since the $\hat{\ell}_0$ approximation is not exact, and might result in some values being closer to zero but not precisely zero. We therefore encourage the attack to focus only on the most influential features, discarding less significant contributions. We then decrease the step size η by following a cosine-annealing schedule (Rony et al., 2018; Pintor et al., 2021), and adjust the sparsity threshold τ dynamically. In particular, if the current sample is adversarial, we increase τ to promote sparser perturbations; otherwise, we decrease τ to reduce \mathcal{L} . The variations of τ are also iteratively reduced following the same cosine-annealing schedule of the step size. The above process is repeated for N iterations, and if during each iteration, we find a better solution that is adversarial and has a lower ℓ_0 norm, we update the optimal perturbation δ^* to the current minimum (line 10). Finally, the best adversarial perturbation δ^* identified during the optimization process is returned (line 11). In conclusion, the main contributions behind σ -zero are: (i) the idea of exploiting the numerically-stable approximation of the ℓ_0 norm by Osborne et al. (2000b) to design a novel loss function (Equation 5), which enables simultaneously searching for an adversarial example while minimizing the ℓ_0 norm of the perturbation (i.e., a non-trivial task given the non-convexity of this norm); and (ii) the introduction of the sparsity threshold τ and its dynamic adjustment policy which, along with gradient normalization and step size annealing, help find very sparse adversarial perturbations faster. The combination of our novel formulation with the aforementioned optimization tricks yields a very fast and reliable ℓ_0 -norm attack algorithm, which does not even require specific hyperparameter tuning, as we will show in our experimental results.

3 EXPERIMENTS

We report the extensive evaluation of the proposed σ -zero attack to compare its performance and efficiency with other state-of-the-art ℓ_0 attacks, considering sixteen baseline and robust models and three different datasets.

3.1 EXPERIMENTAL SETUP

Datasets. We conduct experiments on three popular datasets used for benchmarking adversarial robustness: MNIST (LeCun & Cortes, 2005), CIFAR10 (Krizhevsky, 2009) and ImageNet (Krizhevsky et al., 2012). We use a random subset of 1000 test samples from ImageNet to evaluate attacks performance on it, while we consider the entire test set for MNIST and CIFAR10. For the MNIST and CIFAR10 datasets we used a batch size of 32, while for ImageNet we opted for a batch size of 16.

Attacks. We compare σ -zero against the following state-of-the-art, minimum-norm attacks, in their ℓ_0 -norm variants: the Voting Folded Gaussian Attack (VFGA) attack (Césaire et al., 2021), the Primal-Dual Proximal Gradient Descent (PDPGD) attack (Matyasko & Chau, 2021), the Brendel & Bethge (BB) attack (Brendel et al., 2019), including also its variant with adversarial initialization (BBadv), and the Fast Minimum Norm (FMN) attack (Pintor et al., 2021). We also consider two state-of-the-art ℓ_1 -norm attacks as additional baselines, i.e., the Elastic-Net (EAD) attack (Chen et al., 2018) and SparseFool (Modas et al., 2019), along with two further ℓ_0 -norm attacks, i.e., the ℓ_0 -norm Projected Gradient Descent (PGD- ℓ_0) attack (Croce & Hein, 2019) and the Sparse Random Search (Sparse-RS) attack (Croce et al., 2022).² Compared to minimum-norm attacks, PGD- ℓ_0 and Sparse-RS aim to maximize misclassification confidence within a given maximum number of modifiable features k . Thus, to ensure a fair comparison with minimum-norm attacks, as suggested by Rony et al. (2021b), we tune their perturbation budget k by performing a sample-wise binary search to find minimum-norm adversarial examples. Further details are reported in the Appendix. Finally, we configure all attacks to manipulate input values separately, without constraining the manipulations to individual pixels; e.g., on CIFAR10, the number of modifiable inputs is thus $3 \times 32 \times 32 = 3072$.

Models. We use a selection of both baseline and robust models to evaluate the attacks under different conditions. Our goal is to compare σ -zero on a vast set of models to ensure its broad effectiveness and to expose vulnerabilities that may not be revealed by other attacks (Croce & Hein, 2021a). For the MNIST dataset, we consider two adversarially-trained convolutional neural network (CNN) models by Rony et al. (2021a), i.e., CNN-DDN and CNN-Trades. These models have been trained to be robust to both ℓ_2 and ℓ_∞ adversarial attacks. We denote them respectively with M1 and M2. For the CIFAR10 and ImageNet datasets, we employ state-of-the-art robust models from RobustBench (Croce et al., 2021). For CIFAR10, we adopt eight models, denoted with C1-C10. C1 (Croce et al., 2021) is a non-robust WideResNet-28-10 model. C2 (Carmon et al., 2019) and C3 (Augustin et al., 2020) combine training data augmentation with adversarial training to improve robustness to ℓ_∞ and ℓ_2 attacks. C4 (Engstrom et al., 2019) is an adversarially trained model that is robust to ℓ_2 -norm attacks. C5 (Gowal et al., 2021) exploits generative models to artificially augment the original training set and improve adversarial robustness to generic ℓ_p -norm attacks. C6 (Chen et al., 2020) is a robust ensemble model. C7 (Xu et al., 2023) is a recently proposed adversarial training defense robust to ℓ_2 attacks. C8 (Addepalli et al., 2022) enforces diversity during data augmentation and combines it with adversarial training. Finally, we also include the ℓ_1 robust models C9 (Croce & Hein, 2021b) and C10 (Jiang et al., 2023). For ImageNet, we consider a pretrained ResNet-18 denoted with I1 (He et al., 2015), and five robust models to ℓ_∞ -attacks, denoted with I2 (Engstrom et al., 2019), I3 (Wong et al., 2020), I4 (Salman et al., 2020), I5 (Hendrycks et al., 2021), and I6 (Salman et al., 2020).

Hyperparameters. We conduct our experiments using the default hyperparameters used in the original implementation of the attacks from AdversarialLib (Rony & Ben Ayed) and Foolbox (Rauber et al., 2017). We only change the number of steps to 1000, to ensure that all attacks reach convergence (Pintor et al., 2022). VFGA (Césaire et al., 2021) constitutes the only exception, as it terminates only once an adversarial example is obtained. We report additional results using 100 steps in the Appendix. As gradient-based attacks perform one forward and one backward pass in each step, we double the steps for Sparse-RS, which, being a gradient-free attack, only makes one forward pass per iteration. This ensures a fair comparison. For σ -zero, we set 1000 steps, $\eta_0 = 1$, $\tau_0 = 0.5$ and $\sigma = 0.1$. We keep the same configuration for all models and datasets, showing that no specific hyperparameter tuning is required for σ -zero. Additional analyses of the influence of the hyperparameters on the performance of σ -zero can be found in the Appendix.

Evaluation Metrics. For each attack, we report the Attack Success Rate (ASR), defined as the ratio of successfully attacked samples, and the median ℓ_0 norm. Additionally, we report ASR_k , which indicates the ASR of attacks with a fixed budget of k perturbed features. We also compare the computational effort of each attack considering their execution time, the average number of queries (i.e., the sum of #forwards and #backwards) needed to perform each attack, and the Video Random Access Memory (VRAM) consumption.³ We measure the execution time on a workstation with NVIDIA A100 Tensor Core GPU (40GB memory) and two Intel® Xeon® Gold 6238R processors. For measuring the memory consumption, we consider the maximum amount of VRAM used by each

²Sparse-RS is a gradient-free (black-box) attack, which only requires query access to the target model. We consider it as an additional baseline in our experiments, but it should not be considered a direct competitor of gradient-based attacks, as it works under much stricter assumptions (i.e., no access to input gradients).

³VRAM is a type of memory designed explicitly for use in Graphics Processing Units (GPUs).

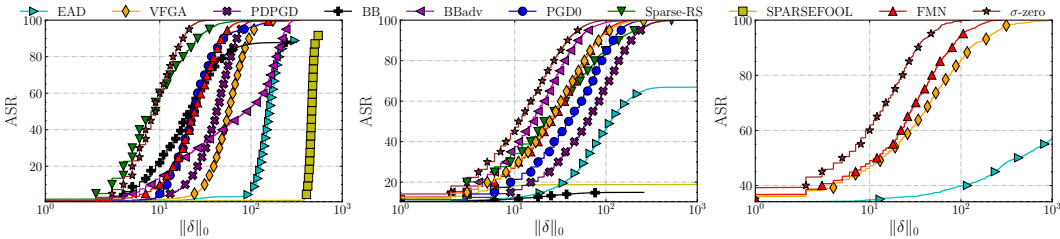


Figure 2: Robustness evaluation curves, reporting ASR versus perturbation size, for M2 on MNIST (leftmost plot), C2 on CIFAR10 (middle plot), and I1 on ImageNet (rightmost plot).

attack among all the batches, which is a minimum requirement to run it without failure. By assessing the performance of each attack across these various metrics, we can gain a more comprehensive understanding of their effectiveness and scalability.

3.2 EXPERIMENTAL RESULTS

Attack Performance. Table 1 reports, for all models and datasets, the median value of $\|\delta\|_0$ and the attack success rates. The values obtained confirm that our attack can find smaller perturbations in all cases. Specifically, over all the dataset-model configurations, $\sigma\text{-zero}$ drastically improves the state of the art of sparse attacks. For example, on CIFAR10 models, $\sigma\text{-zero}$ outperforms FMN by reducing the median number of manipulated features from 52 to 32 in the best case (C9) and from 7 to 5 in the worst case (C1). On ImageNet models, the median $\|\delta\|_0$ is reduced from 58 to 23 in the best case (I6) and from 9 to 3 in the worst case (I2). Furthermore, we observe that the ASR of BB, which is the closest attack in terms of performance to $\sigma\text{-zero}$, drops when used in settings where the input dimensionality increases (e.g., CIFAR10), and it becomes unfeasible in extreme cases (i.e., ImageNet). From Table 1, we can also notice that the median $\|\delta\|_0$ of BB sometimes is ∞ , since its ASR is lower than 50%. BBadv does not suffer from the same issue but $\sigma\text{-zero}$ continues to outperform that variant too. Lastly, we show in the Appendix that our attack always reaches ASR=100% against all models, even when decreasing the number of iterations. For other attacks, this is not ensured, particularly when reducing the number of iterations.

Computational Effort. We report the runtime comparison, the number of queries issued to the model, and the VRAM used by each attack. The results show that our attack is up to 2 (16) times faster than BB when considering MNIST (CIFAR10) models. Therefore, even if BB finds slightly better ℓ_0 -adversarial examples in one configuration, its computational effort is much higher than $\sigma\text{-zero}$. Furthermore, we observed that BB often stops unexpectedly before reaching the specified number of steps because it fails to initialize the attack.

The speed advantage of $\sigma\text{-zero}$ is given because our attack is a simple gradient-based approach that avoids costly inner projections, such as the ones used by BB. On the other hand, $\sigma\text{-zero}$ is slightly slower than FMN and VFGA; however, it compensates by finding better solutions. Notably, similarly to them, $\sigma\text{-zero}$ requires fewer queries than remaining attacks. Furthermore, the speed-competing method VFGA is memory-hungry, forcing us to reduce the batch size when testing its effectiveness on larger models, e.g., C5, C6, and C7. Conversely, running our algorithm also requires reasonable VRAM, as $\sigma\text{-zero}$ implements a lightweight search that includes only the cost of computing gradients and norms for each step. Overall, the practical advantages of our attack make it a promising direction for benchmarking large DNNs in an effective and time-efficient way.

ImageNet Results. For ImageNet, we restrict our analysis to EAD, FMN, and VFGA, as they outperform competing attacks on MNIST and CIFAR10 in terms of ASR, perturbation size, and execution time. While all ImageNet models are deemed robust to ℓ_1 and ℓ_∞ -norm attacks, they are vulnerable to our ℓ_0 -attack. Remarkably, I6 offers higher robustness against ℓ_0 attacks, requiring more effort to evade it. The results show that in most configurations, our attack finds adversarial perturbations with a lower median ℓ_0 -norm, while being at the same time faster and memory-comparable. The results in the Appendix further confirm that even when decreasing the number of iterations to 100, our attack finds lower ℓ_0 -norm solutions and always achieves ASR=100%.

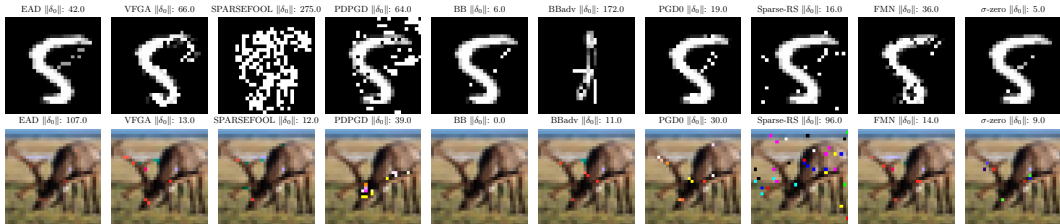


Figure 3: Randomly chosen adversarial examples from MNIST M2 (top-row), CIFAR10 C2 (bottom row) found by adversarial attacks we tested.

Robustness Evaluation Curves. Complementary to the performance results shown in Table 1, we present the robustness evaluation curves in Fig. 2 for each attack on M2, C2, and I1. These curves go beyond the only median statistic and ASR_k , providing more evidence that σ -zero achieves higher ASR with smaller ℓ_0 -norm perturbations compared to the other attacks. Moreover, the ASR of our attack goes up to 100%, validating the correctness of our gradient-based approach even when considering unbounded perturbations (Carlini & Wagner, 2017). These results reinforce our previous findings that σ -zero is an efficient and effective method for generating adversarial examples with smaller ℓ_0 norm. In the Appendix, we include similar curves for all the other experimental configurations, for which results are consistent. In summary, our σ -zero attack consistently outperforms other state-of-the-art methods, suggesting that it can identify smaller and more effective perturbations, making it a highly promising robustness evaluation method.

Visual Inspection of Adversarial Examples. In Fig. 3 we show adversarial examples generated with competing ℓ_0 -attacks, and our σ -zero. First, we can see that ℓ_0 adversarial perturbations are always clearly visually distinguishable. Their goal, indeed, is not to be indistinguishable to the human eye – a common misconception related to adversarial examples (Biggio & Roli, 2018; Gilmer et al., 2018) – but rather to show whether and to what extent models can be fooled by just changing a few input values. For example, note how FMN and VFGA find similar perturbations, as they mostly target overlapping regions of interest. Conversely, EAD finds sparse perturbations scattered throughout the image but with a lower magnitude. This divergence is attributed to EAD’s reliance on an ℓ_1 regularizer, which promotes sparsity, thus diminishing perturbation magnitude without necessarily reducing the number of perturbed features. Conversely, our attack does not focus on specific areas or patterns within the images but identifies diverse critical features, whose manipulation is sufficient to mislead the target models. Given the diversity of solutions that the attacks offer, we argue that their combined usage may still improve adversarial robustness evaluation to sparse attacks.

4 RELATED WORK

Due to the inherent complexity of optimizing over non-convex and non-differentiable constraint, classical gradient-based algorithms like PGD (Madry et al., 2018) cannot be used for computing ℓ_0 -norm attacks. We categorize the existing ℓ_0 -norm attacks into two main groups: (i) multiple-norm attacks extended to ℓ_0 , and (ii) attacks specifically designed to optimize ℓ_0 perturbations. Furthermore, we discuss related work that leverages the approximation of ℓ_0 for different goals.

Multiple-norm attacks extended to ℓ_0 . These attacks are developed to work with multiple ℓ_p norms and include the extension of their algorithms to the ℓ_0 norm. While they are able to find sparse perturbations, they often require strong use of heuristics to work in this setting. Brendel et al. (2019) initializes the attack from an adversarial example far away from the clean sample and optimizes the perturbation by walking with small steps on the decision boundary trying to get closest to the original sample. In general, the algorithm can be used for any ℓ_p norm, including ℓ_0 , but the individual optimization steps are very costly. Pintor et al. (2021) propose the Fast Minimum-Norm (FMN) attack that does not require an initialization step and converges efficiently with lightweight gradient-descent steps. However, their approach was developed to generalize over ℓ_p norms, but it does not make special adaptations to specifically minimize the ℓ_0 norm. Matyasko & Chau (2021) use a two-player approach that optimizes the trade-off between perturbation size and loss of the attack and uses relaxations of the ℓ_0 norm (e.g., $\ell_{1/2}$) to promote sparsity. This scheme however does not strictly minimize the ℓ_0 norm, as the relaxation does not set the lowest components exactly to zero.

ℓ_0 -specific attacks. Croce et al. (2022) introduced SparseRS, a random search-based adversarial attack that explores potential perturbation candidates to return the highest confidence solution. Unlike minimum-norm attacks, their approach is rooted in a maximum-confidence attack framework with a predefined number of feature manipulations. Césaire et al. (2021) have designed an attack specifically for the ℓ_0 norm. This attack is modeled as a stochastic Markov problem. It induces folded Gaussian noise to selected input components, iteratively finding the set that achieves misclassification with minimal perturbation. However, their approach requires a considerable amount of memory to explore the possible combinations and to find an optimal solution. This makes it infeasible to use for larger problems. With σ -zero, we show that the benefits from both groups, efficiency and precision, can be combined to effectively generate sparse ℓ_0 attacks. It stands therefore as a promising solution for evaluating DNNs’ robustness within the ℓ_0 threat model, which remains relatively underexplored in existing benchmarks (Croce et al. 2021).

Approximation of the ℓ_0 norm. Given the nonconvex and discontinuous nature of the ℓ_0 norm, the adoption of surrogate approximation functions has been extensively studied (Bach et al., 2012; Weston et al., 2003; Zhang, 2008). Chen et al. (2018) use elastic-net regularization to calculate sparse perturbations, however, their attack do not necessarily find minimum ℓ_0 -norm perturbations. In our work, we use the formulation proposed by Osborne et al. (2000a), which provides an unbiased estimate of the actual ℓ_0 . Furthermore, it has been employed by Cinà et al. (2022) in the context of poisoning attacks to decrease sparsity in the model’s activations, while we use it as a penalty term for crafting minimum ℓ_0 -norm adversarial examples.

5 CONCLUSIONS, LIMITATIONS, AND FUTURE WORK

Despite numerous proposed attacks for assessing DNN robustness, evaluation methods tend to overlook the significance of ℓ_0 -norm attacks (Chen et al., 2018; Croce & Hein, 2021a). However, these attacks can provide valuable insights into identifying the minimum manipulated input values required for successful attacks and reveal crucial information about model limitations. We argue that this literature gap is primarily due to the non-differentiable nature of the ℓ_0 norm and its computational complexity, which poses challenges for gradient-based optimization.

In this work, we present σ -zero, a novel approach that leverages a smooth approximation of the ℓ_0 norm. By making the objective differentiable, our method becomes amenable to optimization with gradient descent. Through extensive experimentation, we demonstrate the efficacy, precision, and scalability of σ -zero in diverse scenarios, specifically for identifying minimal ℓ_0 perturbations. Our approach consistently discovers smaller minimum-norm perturbations across all models and datasets, while maintaining computational efficiency in execution time and VRAM consumption, and without requiring any computationally-demanding hyperparameter tuning. By identifying the smallest number of input values that can be modified to mislead the target model, our attack provides valuable insights on the vulnerabilities of DNN models and what they learn as salient input characteristics. Additionally, it may also provide meaningful insights on how to mitigate such vulnerabilities to improve robustness.

Although our approach offers promising results for benchmarking DNNs robustness, it relies on the white-box assumption. However, in the absence of such access, attackers may resort to techniques like transferability or gradient estimation to exploit vulnerabilities (Carlini et al., 2019; Tramèr et al., 2020). We acknowledge the significance of this analysis and plan to investigate it further in future research endeavors.

In conclusion, σ -zero emerges as a highly promising candidate for establishing a standardized benchmark to evaluate robustness against sparse ℓ_0 perturbations. By facilitating more reliable and scalable assessments, it is poised to drive significant advancements in the development of novel models with improved robustness guarantees against the specific threat model under consideration.

Ethics Statement. Based on our comprehensive analysis, we assert that there are no identifiable ethical considerations or foreseeable negative societal consequences that warrant specific attention within the confines of this study. Rather this study will help improve the understanding of adversarial robustness properties of DNNs, and identify potential ways in which robustness can be improved.

REFERENCES

- Pravanti Addepalli, Samyak Jain, and Venkatesh Babu R. Efficient and effective augmentation strategy for adversarial training. In *NeurIPS*, 2022.
- Maximilian Augustin, Alexander Meinke, and Matthias Hein. Adversarial robustness on in- and out-distribution improves explainability. In *Computer Vision - ECCV 2020 - 16th European Conference*, volume 12371 of *Lecture Notes in Computer Science*, pp. 228–245. Springer, 2020.
- Francis R. Bach, Rodolphe Jenatton, Julien Mairal, and Guillaume Obozinski. Optimization with sparsity-inducing penalties. *Found. Trends Mach. Learn.*, 4(1):1–106, 2012.
- Battista Biggio and Fabio Roli. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, 84:317–331, 2018.
- Battista Biggio, Iginio Corona, Davide Maiorca, Blaine Nelson, Nedim Srndic, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD*, volume 8190 of *Lecture Notes in Computer Science*, pp. 387–402. Springer, 2013.
- Wieland Brendel, Jonas Rauber, Matthias Kümmerer, Ivan Ustyuzhaninov, and Matthias Bethge. Accurate, reliable and fast robustness evaluation. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems, NeurIPS*, 2019.
- Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy SP*, pp. 39–57. IEEE Computer Society, 2017.
- Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian J. Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness. *CoRR*, abs/1902.06705, 2019.
- Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, John C Duchi, and Percy S Liang. Unlabeled data improves adversarial robustness. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- Manon Césaire, Lucas Schott, Hatem Hajri, Sylvain Lamprier, and Patrick Gallinari. Stochastic sparse adversarial attacks. In *33rd IEEE International Conference on Tools with Artificial Intelligence, ICTAI*, pp. 1247–1254. IEEE, 2021.
- Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. EAD: elastic-net attacks to deep neural networks via adversarial examples. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18)*, pp. 10–17. AAAI Press, 2018.
- Tianlong Chen, Sijia Liu, Shiyu Chang, Yu Cheng, Lisa Amini, and Zhangyang Wang. Adversarial robustness: From self-supervised pre-training to fine-tuning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, pp. 696–705. Computer Vision Foundation / IEEE, 2020.
- Antonio Emanuele Cinà, Ambra Demontis, Battista Biggio, Fabio Roli, and Marcello Pelillo. Energy-latency attacks via sponge poisoning. *CoRR*, abs/2203.08147, 2022.
- Francesco Croce and Matthias Hein. Sparse and imperceptible adversarial attacks. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4723–4731, 2019.
- Francesco Croce and Matthias Hein. Mind the box: l_1 -apgd for sparse adversarial attacks on image classifiers. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML*, volume 139 of *Proceedings of Machine Learning Research*, pp. 2201–2211. PMLR, 2021a.
- Francesco Croce and Matthias Hein. Mind the box: l_1 -apgd for sparse adversarial attacks on image classifiers. In *International Conference on Machine Learning (ICML)*, 2021b.

- Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo DeBenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks*, 2021.
- Francesco Croce, Maksym Andriushchenko, Naman D. Singh, Nicolas Flammarion, and Matthias Hein. Sparse-rs: A versatile framework for query-efficient sparse black-box adversarial attacks. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI*, pp. 6437–6445. AAAI Press, 2022.
- Edoardo DeBenedetti, Vikash Sehwal, and Prateek Mittal. A light recipe to train robust vision transformers. In *First IEEE Conference on Secure and Trustworthy Machine Learning*, 2023. URL <https://openreview.net/forum?id=IztT98ky0cKs>.
- Logan Engstrom, Andrew Ilyas, Hadi Salman, Shibani Santurkar, and Dimitris Tsipras. Robustness (python library), 2019. URL <https://github.com/MadryLab/robustness>
- Justin Gilmer, Ryan P. Adams, Ian J. Goodfellow, David Andersen, and George E. Dahl. Motivating the rules of the game for adversarial example research. *CoRR*, abs/1807.06732, 2018.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations (ICLR)*, 2015.
- Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and Timothy A. Mann. Improving robustness using generated data. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS*, pp. 4218–4233, 2021.
- Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.
- Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV*, pp. 8320–8329. IEEE, 2021.
- Yulun Jiang, Chen Liu, Zhichao Huang, Mathieu Salzmann, and Sabine Süsstrunk. Towards stable and efficient adversarial training against l1 bounded adversarial attacks. In *International Conference on Machine Learning*, 2023.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60:84 – 90, 2012.
- Yann LeCun and Corinna Cortes. The mnist database of handwritten digits. 2005.
- A Madry, A Makelov, L Schmidt, D Tsipras, and A Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2018.
- Alexander Matyasko and Lap-Pui Chau. PDPGD: primal-dual proximal gradient descent adversarial attack. *CoRR*, abs/2106.01538, 2021. URL <https://arxiv.org/abs/2106.01538>.
- Apostolos Modas, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Sparsefool: a few pixels make a big difference. In *Conference on computer vision and pattern recognition (CVPR)*, 2019.
- Michael R. Osborne, Brett Presnell, and Berwin A. Turlach. On the lasso and its dual. *Journal of Computational and Graphical Statistics*, 9:319 – 337, 2000a.
- Michael R. Osborne, Brett Presnell, and Berwin A. Turlach. On the lasso and its dual. *Journal of Computational and Graphical Statistics*, 9:319–337, 2000b.
- Maura Pintor, Fabio Roli, Wieland Brendel, and Battista Biggio. Fast minimum-norm adversarial attacks through adaptive norm constraints. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems, NeurIPS*, pp. 20052–20062, 2021.

- Maura Pintor, Luca Demetrio, Angelo Sotgiu, Ambra Demontis, Nicholas Carlini, Battista Biggio, and Fabio Roli. Indicators of attack failure: Debugging and improving optimization of adversarial examples. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 23063–23076. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/91ffdc5e2f12436d99914418e38d0a09-Paper-Conference.pdf.
- Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models, 2017. URL <https://github.com/bethgelab/foolbox>.
- Jérôme Rony, Luiz G. Hafemann, Luiz Oliveira, Ismail Ben Ayed, Robert Sabourin, and Eric Granger. Decoupling direction and norm for efficient gradient-based l2 adversarial attacks and defenses. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4317–4325, 2018.
- Jérôme Rony, Eric Granger, Marco Pedersoli, and Ismail Ben Ayed. Augmented lagrangian adversarial attacks. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV*, pp. 7718–7727. IEEE, 2021a.
- Jérôme Rony, Eric Granger, Marco Pedersoli, and Ismail Ben Ayed. Augmented lagrangian adversarial attacks. In *Conference on computer vision and pattern recognition (CVPR)*, 2021b.
- Jérôme Rony and Ismail Ben Ayed. Adversarial Library. URL <https://github.com/jeromerony/adversarial-library>.
- Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS*, 2020.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations (ICLR)*, 2014a.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Yann LeCun. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR*, 2014b.
- Florian Tramèr, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS*, 2020.
- Jason Weston, André Elisseeff, Bernhard Schölkopf, and Michael E. Tipping. Use of the zero-norm with linear models and kernel methods. *J. Mach. Learn. Res.*, 3:1439–1461, 2003.
- Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training. In *8th International Conference on Learning Representations, ICLR*. OpenReview.net, 2020.
- Yuancheng Xu, Yanchao Sun, Micah Goldblum, Tom Goldstein, and Furong Huang. Exploring and exploiting decision boundary dynamics for adversarial robustness. In *International Conference on Learning Representations (ICLR)*, 2023.
- Tong Zhang. Multi-stage convex relaxation for learning with sparse regularization. In *NIPS*, 2008.