

Formulation

Consider a l_0 bounded adversarial perturbation δ , we can decompose it into a magnitude tensor \mathbf{p} and a binary sparsity mask \mathbf{m} . Therefore, the attacker aims to maximize the following objective function:

$$\max_{\|\delta\|_0 \leq k, 0 \leq \mathbf{x} + \delta \leq 1} \mathcal{L}(\theta, \mathbf{x} + \delta) = \max_{\mathbf{p} \in \mathcal{S}_p, \mathbf{m} \in \mathcal{S}_m} \mathcal{L}(\theta, \mathbf{x} + \mathbf{p} \odot \mathbf{m})$$

The feasible sets for \mathbf{p} and \mathbf{m} are $\mathcal{S}_p = \{\mathbf{p} \in \mathbb{R}^{h \times w \times c} | 0 \leq \mathbf{x} + \mathbf{p} \leq 1\}$ and $\mathcal{S}_m = \{\mathbf{m} \in \{0, 1\}^{h \times w \times 1} | \|\mathbf{m}\|_0 \leq k\}$, respectively. k is the sparsity level.

Sparse-PGD (sPGD)

Similar to PGD, sPGD iteratively updates \mathbf{p} and \mathbf{m} until finding a successful adversarial example or reaching the maximum iteration number.

Update Magnitude Tensor \mathbf{p} : We utilize PGD in the l_∞ case, i.e., use the sign of the gradients, to optimize \mathbf{p} , with α being the step size:

$$\mathbf{p} \leftarrow \Pi_{\mathcal{S}_p}(\mathbf{p} + \alpha \cdot \text{sign}(\nabla_{\mathbf{p}} \mathcal{L}(\theta, \mathbf{x} + \mathbf{p} \odot \mathbf{m})))$$

Update Sparsity Mask \mathbf{m} : Instead of optimizing the discrete \mathbf{m} , we update its continuous alternative $\tilde{\mathbf{m}} \in \mathbb{R}^{h \times w \times 1}$, and then project $\tilde{\mathbf{m}}$ to the feasible set \mathcal{S}_m to obtain \mathbf{m} :

$$\begin{aligned} \tilde{\mathbf{m}} &\leftarrow \tilde{\mathbf{m}} + \beta \cdot \nabla_{\tilde{\mathbf{m}}} \mathcal{L} / \|\nabla_{\tilde{\mathbf{m}}} \mathcal{L}\|_2 \\ \mathbf{m} &\leftarrow \Pi_{\mathcal{S}_m}(\sigma(\tilde{\mathbf{m}})) \end{aligned}$$

where β is the step size, $\sigma(\cdot)$ denotes sigmoid function. $\Pi_{\mathcal{S}_m}$ is to set the k -largest elements in $\tilde{\mathbf{m}}$ to 1 and the rest to 0.

Backward Function: We propose the *sparse/projected* gradient g_p and *unprojected* gradient \tilde{g}_p of \mathbf{p} , which exhibit complementary performance:

$$\begin{aligned} g_p &= \nabla_{\mathbf{p}} \mathcal{L}(\theta, \mathbf{x} + \delta) \odot \mathbf{m} \\ \tilde{g}_p &= \nabla_{\mathbf{p}} \mathcal{L}(\theta, \mathbf{x} + \delta) \odot \sigma(\tilde{\mathbf{m}}) \end{aligned}$$

Random Reinitialization: When the attack fails and the current sparsity mask \mathbf{m} remains unchanged for 3 consecutive iterations, $\tilde{\mathbf{m}}$ will be randomly reinitialized.

Pseudo-code

The pseudo-code of sPGD is presented below:

Algorithm 1 Sparse-PGD

```

1: Input: Clean image:  $\mathbf{x} \in [0, 1]^{h \times w \times c}$ ; Model parameters:  $\theta$ ; Max iteration number:  $T$ ; Tolerance:  $t$ ;  $l_0$  budget:  $k$ ; Step size:  $\alpha, \beta$ ; Small constant:  $\gamma = 2 \times 10^{-8}$ 
2: Random initialize  $\mathbf{p}$  and  $\tilde{\mathbf{m}}$ 
3: for  $i = 0, 1, \dots, T - 1$  do
4:    $\mathbf{m} = \Pi_{\mathcal{S}_m}(\sigma(\tilde{\mathbf{m}}))$ 
5:   Calculate the loss  $\mathcal{L}(\theta, \mathbf{x} + \mathbf{p} \odot \mathbf{m})$ 
6:   if unprojected then
7:      $g_p = \nabla_{\mathbf{p}} \mathcal{L} \odot \sigma(\tilde{\mathbf{m}})$             $\{\delta = \mathbf{p} \odot \mathbf{m}\}$ 
8:   else
9:      $g_p = \nabla_{\mathbf{p}} \mathcal{L} \odot \mathbf{m}$ 
10:  end if
11:   $\tilde{g}_p = \nabla_{\mathbf{p}} \mathcal{L} \odot \mathbf{p} \odot \sigma'(\tilde{\mathbf{m}})$ 
12:   $\mathbf{p} = \Pi_{\mathcal{S}_p}(\mathbf{p} + \alpha \cdot \text{sign}(g_p))$ 
13:   $\mathbf{d} = \tilde{g}_p / (\|\tilde{g}_p\|_2)$  if  $\|\tilde{g}_p\|_2 \geq \gamma$  else 0
14:   $\mathbf{m}_{old}, \tilde{\mathbf{m}} = \mathbf{m}, \tilde{\mathbf{m}} + \beta \cdot \mathbf{d}$ 
15:  if attack succeeds then
16:    break
17:  end if
18:  if  $\|\Pi_{\mathcal{S}_m}(\sigma(\tilde{\mathbf{m}})) - \mathbf{m}_{old}\|_0 \leq 0$  for  $t$  consecutive iters then
19:    Random initialize  $\tilde{\mathbf{m}}$ 
20:  end if
21: end for
22: Output: Perturbation:  $\delta = \mathbf{p} \odot \mathbf{m}$ 

```

Sparse-AutoAttack (sAA)

For comprehensive robustness evaluation against l_0 bounded perturbations, we propose **sparse-AutoAttack (sAA)**, which is a parameter-free ensemble of both white-box and black-box attacks. We adopt two variants of sPGD with different backward functions: **sPGD_{proj}** using g_p , and **sPGD_{unproj}** using \tilde{g}_p . Additionally, **Sparse-RS** is adopted as the black-box attack. These attacks run in a cascade way.

Adversarial Training

To accommodate the scenario of adversarial training, we make the following modifications to sPGD: **a)** random backward function, **b)** multi- k strategy, **c)** higher tolerance for reinitialization. In experiments, we incorporate sPGD in vanilla adversarial training and TRADES and name corresponding methods **sAT** and **sTRADES**, respectively.

Experiments

Table 1. Robust accuracy of various models on different attacks that generate l_0 bounded perturbations, where the sparsity level $k = 20$. The models are trained on **CIFAR-10**.

Model	Network	Clean	Black-Box			White-Box			sAA	
			CS	RS	SF	PGD ₀	SAIF	sPGD _{proj}		sPGD _{unproj}
Vanilla	RN-18	93.9	1.2	0.0	17.5	0.4	3.2	0.0	0.0	0.0
<i>l_∞-adv. trained, $\epsilon = 8/255$</i>										
GD	PRN-18	87.4	26.7	6.1	52.6	25.2	40.4	9.0	15.6	5.3
PORT	RN-18	84.6	27.8	8.5	54.5	21.4	42.7	9.1	14.6	6.7
<i>l₁-adv. trained, $\epsilon = 12$</i>										
l ₁ -APGD	PRN-18	80.7	32.3	25.0	65.4	39.8	55.6	17.9	18.8	16.9
Fast-EG-l ₁	PRN-18	76.2	35.0	24.6	60.8	37.1	50.0	18.1	18.6	16.8
<i>l₀-adv. trained, $k = 20$</i>										
PGD ₀ -A	PRN-18	77.5	16.5	2.9	62.8	56.0	47.9	9.9	21.6	2.4
PGD ₀ -T	PRN-18	90.0	24.1	4.9	85.1	61.1	67.9	27.3	37.9	4.5
sAT	PRN-18	84.5	52.1	36.2	81.2	78.0	76.6	75.9	75.3	36.2
sTRADES	PRN-18	89.8	69.9	61.8	88.3	86.1	84.9	84.6	81.7	61.7

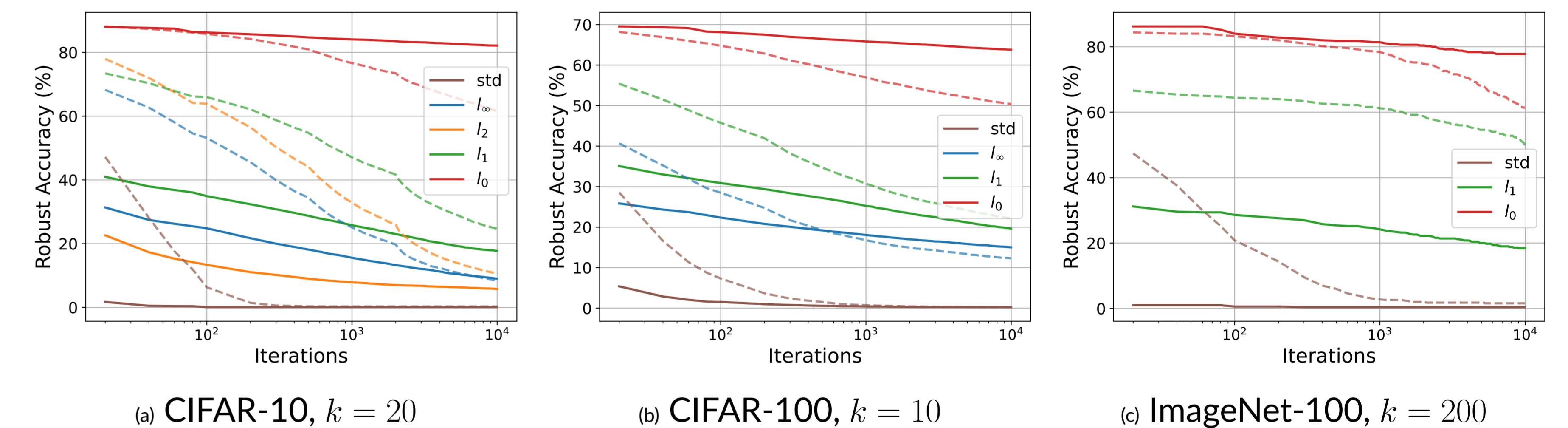


Figure 1. Comparison between sPGD and RS attack under different iterations. The results of sPGD and RS attack are shown in solid lines and dotted lines, respectively.

Takeaway Messages

1. We propose **sparse-PGD (sPGD)** to generate l_0 bounded adversarial perturbations.
2. We combine sPGD with a black-box attack to construct **sparse-AutoAttack (sAA)** for a more comprehensive robustness evaluation against sparse attacks.
3. sPGD achieves better performance under limited iterations. Models adversarially trained by sPGD have the strongest robustness against sparse attacks.

Codes on Github

