

Problem Formulation

Generating adversarial images that fool a neural network can be formalized as an optimization problem with constraints.

Let f be a classification function that takes as input an image sample, and produces a vector containing c containing the distribution of probability for each class.

Given an image x , and the associated class y , that is classified correctly by the function f , the goal of an attack is to change the prediction of the function, by producing an adversarial image \bar{x} such that $f(\bar{x}) \neq y$. The image \bar{x} must be similar to x , in such a way that the L_p norm of the difference is bounded by a constant ϵ :

$$\arg \max_i f_i(\bar{x}) \neq y \quad \|x - \bar{x}\|_p \leq \epsilon \quad (1)$$

with $\epsilon \in \mathbb{N}_+$ in our black-box case. An attack can be untargeted, as described above, or targeted, by forcing the miss-classification of the adversarial image to a specific adversarial class $\bar{y} \neq y$: $f(\bar{x}) = \bar{y}$. The task of finding the perturbed image \bar{x} , associated to x , can be viewed as a minimization problem

$$\min_{\bar{x}} L(f(\bar{x}), y) \quad \|x - \bar{x}\|_p \leq \epsilon, \quad (2)$$

whith L equals to $L(f(\bar{x}), y) = f_y(\bar{x})$ for untargeted attacks and $L(f(\bar{x}), \bar{y}) = 1 - f_{\bar{y}}(\bar{x})$ for targeted ones.

Pixle

Pixle is a black-box attack based on random search. Despite its simplicity, random search performs well in many situations and does not depend on gradient information associated to the objective function $L(\cdot, \cdot)$.

It attacks an image x , having class y , by moving the pixels in it. The algorithm performs R restarts, and T iterations per each restart. At each iteration, the algorithm randomly select a random source patch from the image and, for each pixel in it, it calculates a destination pixel using a function m . If the loss decreases, then the replacement is kept in the next iteration. At the end of each iteration the attack that reduced the loss more is permanently added to the adversarial image. Pixle is part of TorchAttacks framework (Kim 2020).

Pseudocode

Algorithm 1 The overall attacking procedure

Require: input image x with its associated label y . Maximum and minimum dimension for source patch. The number of restarts R and the iterations to perform for each restart step T . The mapping function m .

```

1:  $\bar{x} \leftarrow x$ 
2:  $l \leftarrow f_y(x)$ 
3: for  $r = 0$  to  $R$  do
4:    $x^r \leftarrow \bar{x}$ 
5:   for  $t = 0$  to  $T$  do
6:     Sample  $p = (o_x, o_y, w_p, h_p)$ .
7:     Calculate the set  $P$ 
8:      $x^t \leftarrow \bar{x}$ 
9:     for  $\forall (i, j) \in P$  do
10:       $(z, k) \leftarrow m(i, j)$ 
11:       $x_{z,k}^t \leftarrow x_{i,j}$ 
12:    end for
13:    if  $L(f(x^t), y) < l$  then
14:       $l \leftarrow L(f(x^t), y)$ 
15:       $x^r \leftarrow x^t$ 
16:    end if
17:  end for
18:   $\bar{x} \leftarrow x^r$ 
19: end for
20: return  $\bar{x}$ 

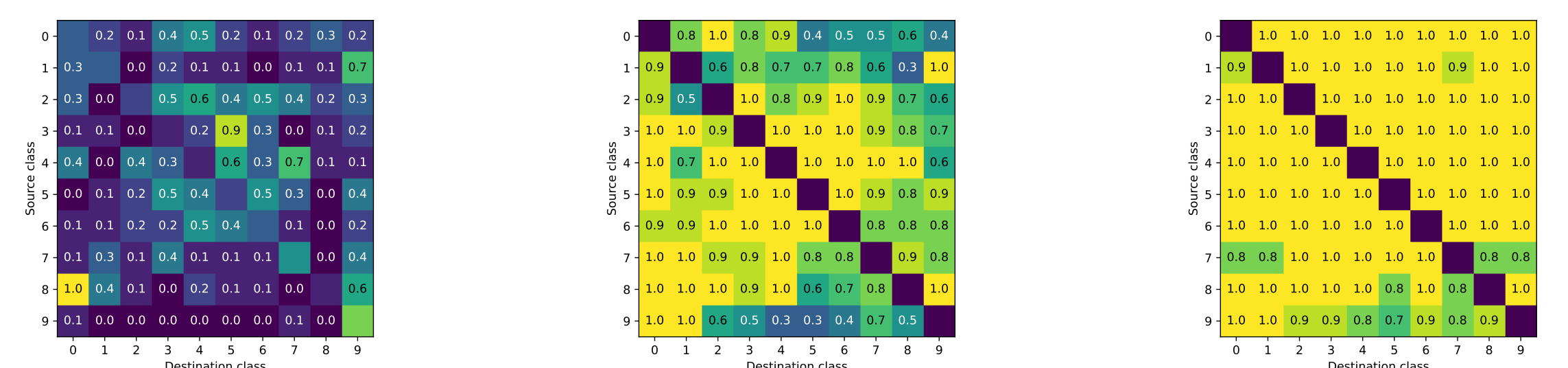
```

Results: success rate, iterations, and L_0 norm

Table 1: Results obtained when attacking multiple datasets trained on ResNet and VGG architectures. For each score, we show the mean and the variance, if present, calculated over all the images on that dataset. Best results for each pair dataset-architecture are highlighted in bold. We compared our proposal to OnePixel Su et al. 2019 and Scratch That Jere et al. 2019

Dataset	Model	Method	Success rate	Iterations	L_0 norm
CIFAR10	ResNet18	OnePixel	64.7	5125 \pm 799	5
		ScratchThat	99.7	1010	38.3 \pm 5.2
		Pixle (Proposed)	100	119 \pm 141	26.8 \pm 22.8
	VGG11	OnePixel	84.5	5100	5
		ScratchThat	99.3	1010	27.64 \pm 4.8
		Pixle (Proposed)	100	80 \pm 145	20.1 \pm 21.9
TinyImageNet	ResNet50	OnePixel	21.0	5100	5
		ScratchThat	69.7	1010	49.1 \pm 2.2
		Pixle (Proposed)	99.6	310 \pm 561	59.0 \pm 88.5
	VGG16	OnePixel	31.9	5100	5
		ScratchThat	76.6	1010	48.6 \pm 2.4
		Pixle (Proposed)	100	87 \pm 201	21.5 \pm 30.6
ImageNet	ResNet50	OnePixel	47.7	5100	5
		ScratchThat	82.6	623 \pm 321	175.2 \pm 25.7
		Pixle (Proposed)	98.0	341 \pm 426	155.7 \pm 184.2
	VGG16	OnePixel	31.9	5100	5
		ScratchThat	81.8	753 \pm 156	143.0 \pm 6.0
		Pixle (Proposed)	99.0	519 \pm 780	98.5 \pm 137.5

Results: targeted attack



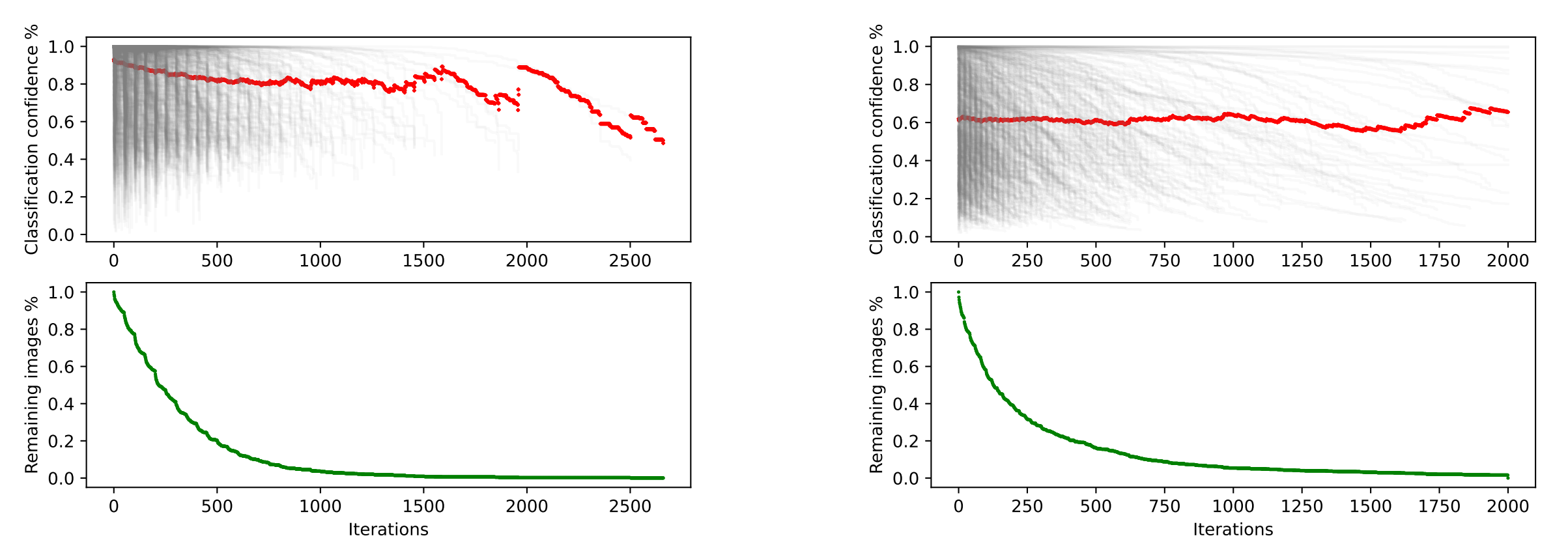
(a) The results obtained using OnePixel attack.

(b) The results obtained using Scratch That attack

(c) The results obtained using our proposal.

Figure 2: The images show the success rate on ResNet18 trained using CIFAR10. The results are calculated using 20 test images, classified correctly by the model, for each class. Each matrix contains the percentage of attacks that have been successfully completed.

Results: iterations and loss



(a) CIFAR10 results using VGG11.

(b) ImageNet results using ResNet50.

Figure 3: Each figure shows, on top, how the losses (the probability associated to the correct class) change during the iterations of our proposal, using the Restart-Iterative algorithm (the red dots are the average loss calculated on that iteration); while the bottom image shows how many images are left to attack after each iteration. Better viewed in colors.



References

- Kim, H. (2020). "Torchattacks: A pytorch repository for adversarial attacks". In: *arXiv preprint arXiv:2010.01950*.
- Jere, M., L. Rossi, B. Hitaj, G. Ciocarlie, G. Boracchi, and F. Koushanfar (2019). "Scratch that! An evolution-based adversarial attack against neural networks". In: *arXiv preprint arXiv:1912.02316*.
- Su, J., D. V. Vargas, and K. Sakurai (2019). "One pixel attack for fooling deep neural