

Adversarial Attacks

Rassin, Royi

June 30, 2021

1 Link to our work

<https://colab.research.google.com/drive/1LyLVopLiKEvfT1VzQ3d1bsiKsh9EG8qN?usp=sharing>

2 The assignment

In this work we are tasked to attack a Resnet18 using a Whitebox and a Black-box approach, while training on the SVHN dataset. We attack with Projected Gradient Descent, and conduct some experiments to explore the optimal hyper-parameters for attack and defense.

3 Training on the SVHN dataset

Using a pretrained ResNet18 means the training process was relatively simple, and all we had to do is attach a new head; a linear layer with an output-dim of 10, for the SVHN dataset. In terms of hyper-parameters for the ResNet18, we used an Adam Optimizer with default parameters, CrossEntropy with a 'sum' reduction, and a batch-size of 64.

4 WhiteBox

4.1 No Defense

In this setting, we simply train the network to solve the SVHN dataset for five epochs, and then apply the attack on the test-set. The parameters of the PGD's $\epsilon = 0.03$, the step-size was set to 0.008, and the number of steps to four. As expected, the accuracy on the perturbed dataset was abysmal, while the natural version was high (you can see the accuracies of all models on table-1, below).

4.2 With Defense

Here, for every sample from the 'natural' dataset, we also train on its perturbed version. The well-known trade-off of defending our networks from attacks is

between accuracy and robustness, and as such, we improved the accuracy on the perturbed dataset to 54%, however, the accuracy on the natural dataset dropped to 91.4%. Further training the model on the perturbed dataset would probably lead to higher accuracy on it and lower on the natural version, but the increase of accuracy every epoch was not significant (from 51 to 52...).

5 Blackbox

5.1 Baseline to Defended

We fed the network that was trained on perturbed data the adversarial samples that were generated from the baseline model, and it fared much better than the baseline, achieving 83.9 accuracy.

5.2 Defended to Baseline

This time, the accuracy dropped significantly, to 55.

5.3 Tinkering with ϵ

After changing the number of steps in the PGD attack to two, we began testing out different epsilons. A greater epsilon will result in a bolder permutation of the image, and a lower chance for the attack to succeed, while a lesser one, will not change the data by much, and as a result, may reduce the accuracy by more (but the change to the image is less apparent). We tinkered with epsilon on the baseline model and the defended model. The baseline's accuracy dropped to 4% while the defended to 31.3%. Graphs and images of the tinkered images can be seen in the appendix. As previously mentioned, the more apparent the difference between the original image to the tinkered one, the greater the epsilon. In the appendix you can see the graphs of epsilon over accuracy per model, and some of the tinkered images (there was no apparent difference among the images from the baseline or the defended one). Regarding the question of which epsilons resulted in a successful attack and which ones did not, we find that the lower the epsilon, the more successful the attack. The reason for that is that the image still resembles the original, yet the accuracy drops significantly (to 25 for the baseline, and 70 for the defended). However, the most perturbed images are not as effective, since they are quite different from the original. Nevertheless, the images I presented here all still contain the digit '3', and I would expect a network to classify it as such, and if it does not (as we can clearly see from increasing the epsilon), the attack indeed works. So all in all, as long as the image still contains the same content (there is a '3' in the image), and the network was fooled to think otherwise, then it worked.

Attack Setting	Defense	Accuracy
None	None	92.9
None	Y	91.4
Whitebox	None	7.3
Whitebox	Y	54
Blackbox	Y	83.9
Blackbox	None	55

6 Appendix

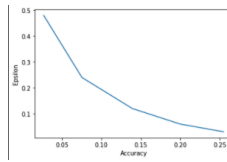


Figure 1: Baseline Epsilons

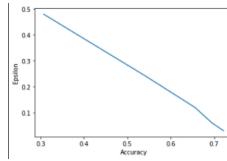


Figure 2: Defended Epsilons

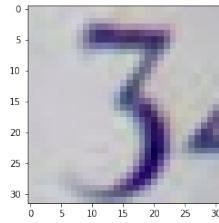


Figure 3: Original Image

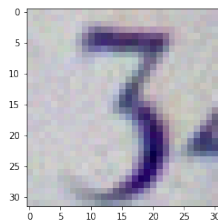


Figure 4: Least tinkered Image

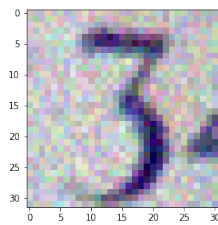


Figure 5: Tinkered Image

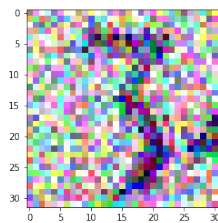


Figure 6: Most tinkered Image