
Supplementary Material

Yanzhi Chen¹ Renjie Xie² Zhanxing Zhu³

A. Theory: Proof of Theorem

Theorem 1. Let T' be the number of backtracking steps. Assuming that (a) the solved \mathbf{z}^* is optimal in each \mathbf{z} -step (b) the backtracking learning rate η' satisfies $\eta' \|\mathbf{H}\| \leq 1$ where $\mathbf{H} = \nabla_{\mathbf{z}^*, \mathbf{z}^*}^2 l_G$ is the Hessian of l_G at $(\mathbf{z}^*, \mathbf{x}')$ and (c) \mathbf{z}_{new}^0 is not far away from \mathbf{z}^* and l_G is smooth w.r.t \mathbf{z} . Then $\nabla_{\mathbf{x}'} Q(\mathbf{x}')$ is approximately a T' -order approximation to $\nabla_{\mathbf{x}'} \mathbf{z}^*$:

$$\nabla_{\mathbf{x}'} \mathbf{z}^* = \sum_{t=0}^{\infty} \mathbf{B} \mathbf{A}^t, \quad \nabla_{\mathbf{x}'} Q(\mathbf{x}') \approx \sum_{t=0}^{T'} \mathbf{B} \mathbf{A}^t$$

where \mathbf{A} and \mathbf{B} is some matrix with $\|\mathbf{A}\| \leq 1$.

Proof: We first prove $\nabla_{\mathbf{x}'} \mathbf{z}^* = \sum_{t=0}^{\infty} \mathbf{B} \mathbf{A}^t$, then $\nabla_{\mathbf{x}'} Q(\mathbf{x}') \approx \sum_{t=0}^{T'} \mathbf{B} \mathbf{A}^t$. This proof is inspired by (Shaban et al., 2019).

Lemma 1. Given conditions (a)(b) in Theorem 1, we have

$$\nabla_{\mathbf{x}'} \mathbf{z}^* = \sum_{t=0}^{\infty} \mathbf{B} \mathbf{A}^t$$

where

$$\mathbf{A} = \mathbf{I} - \eta' \nabla_{\mathbf{z}^*, \mathbf{z}^*}^2 l_G(\mathbf{x}', \mathbf{z}^*), \quad \mathbf{B} = -\eta' \nabla_{\mathbf{x}', \mathbf{z}^*} l_G(\mathbf{x}', \mathbf{z}^*)$$

Proof: it can be easily shown by the technique of implicit differentiation that (Gould et al., 2016):

$$\nabla_{\mathbf{x}'} \mathbf{z}^* = -\nabla_{\mathbf{x}', \mathbf{z}^*} l_G(\mathbf{x}', \mathbf{z}^*) \underbrace{\nabla_{\mathbf{z}^*, \mathbf{z}^*}^2 l_G(\mathbf{x}', \mathbf{z}^*)^{-1}}_{\mathbf{H}^{-1}} \quad (1)$$

Remark that for a matrix \mathbf{M} that $\|\mathbf{M}\| \leq 1$, we have (Horn & Johnson, 2012):

$$(\mathbf{I} - \mathbf{M})^{-1} = \sum_{t=0}^{\infty} \mathbf{M}^t \quad (2)$$

then, since $\eta' \|\mathbf{H}\| \leq 1$,

$$\mathbf{H}^{-1} = \eta' \cdot (\mathbf{I} - (\mathbf{I} - \eta' \mathbf{H}))^{-1} = \eta' \sum_{t=0}^{\infty} (\mathbf{I} - \eta' \mathbf{H})^t \quad (3)$$

substituting this back to (1):

$$\nabla_{\mathbf{x}'} \mathbf{z}^* = \underbrace{-\eta' \nabla_{\mathbf{x}', \mathbf{z}^*} l_G(\mathbf{x}', \mathbf{z}^*)}_{\mathbf{B}} \sum_{t=0}^{\infty} \underbrace{(\mathbf{I} - \eta' \mathbf{H})^t}_{\mathbf{A}} \quad (4)$$

which completes the proof. □

^{*}Equal contribution ¹School of Informatics, The University of Edinburgh, UK ²School of Information Engineering, Southeast University, China ³School of Mathematical Sciences, Peking University, China. Correspondence to: Zhanxing Zhu <zhanxing.zhu@pku.edu.cn>.

Lemma 2. Given conditions (a)(c) in Theorem 1, we have

$$\nabla_{\mathbf{x}'} Q(\mathbf{x}') \approx \sum_{t=0}^{T'} \mathbf{B} \mathbf{A}^t$$

where \mathbf{A} and \mathbf{B} are the same matrices as in Lemma 1.

Proof: remark that

$$\begin{aligned} Q(\mathbf{x}') &:= Q_{\mathbf{x}'}(\mathbf{z}^0) = \underbrace{q_{\mathbf{x}'} \circ q_{\mathbf{x}'} \dots \circ q_{\mathbf{x}'}(\mathbf{z}_0)}_{T'} \\ q_{\mathbf{x}'}(\mathbf{z}) &= \mathbf{z} - \eta' \cdot \nabla_{\mathbf{z}} l_G(\mathbf{z}; \mathbf{x}') \end{aligned} \quad (5)$$

For clarity let us write $q_{\mathbf{x}'}(\mathbf{z}) = q(\mathbf{z}, \mathbf{x}')$ from now on. The total derivative w.r.t \mathbf{x}' in $q(\mathbf{z}, \mathbf{x}')$ is then

$$\nabla_{\mathbf{x}'} q(\mathbf{z}, \mathbf{x}') = \nabla_{\mathbf{z}} q(\mathbf{z}, \mathbf{x}') \Big|_{\mathbf{x}'} \nabla_{\mathbf{x}'} \mathbf{z} + \nabla_{\mathbf{x}'} q(\mathbf{z}, \mathbf{x}') \Big|_{\mathbf{z}} \quad (6)$$

Remark that $\mathbf{z}^{t+1} = q_{\mathbf{x}'}(\mathbf{z}^t) = q(\mathbf{z}^t, \mathbf{x}')$. Then

$$\nabla_{\mathbf{x}'} \mathbf{z}^{t+1} = \underbrace{\nabla_{\mathbf{z}^t} \mathbf{z}^{t+1} \Big|_{\mathbf{x}'}}_{\alpha_t} \nabla_{\mathbf{x}'} \mathbf{z}^t + \underbrace{\nabla_{\mathbf{x}'} \mathbf{z}^{t+1} \Big|_{\mathbf{z}^t}}_{\beta_t} \quad (7)$$

Unfolding this recursive expression, we obtain

$$\nabla_{\mathbf{x}'} \mathbf{z}^{T'} = \sum_{t=0}^{T'} \beta_t \alpha_{t+1} \alpha_{t+2} \dots \alpha_{T'} \quad (8)$$

Now, since $\mathbf{z}^{T'} = \mathbf{z}^*$ (condition (a)), we have

$$\begin{aligned} \alpha_{T'} &= \nabla_{\mathbf{z}^{T'}} \mathbf{z}^{T'+1} \Big|_{\mathbf{x}'} = \nabla_{\mathbf{z}^{T'}} \mathbf{z}^{T'} - \eta' \nabla_{\mathbf{z}^{T'}} \nabla_{\mathbf{z}^{T'}} l_G(\mathbf{z}^{T'}; \mathbf{x}') = \mathbf{I} - \eta' \nabla_{\mathbf{z}^*, \mathbf{z}^*}^2 l_G(\mathbf{x}', \mathbf{z}^*) = \mathbf{A} \\ \beta_{T'} &= \nabla_{\mathbf{x}'} \mathbf{z}^{T'+1} \Big|_{\mathbf{z}^{T'}} = \mathbf{0} - \eta' \nabla_{\mathbf{x}'} \nabla_{\mathbf{z}^{T'}} l_G(\mathbf{z}^{T'}; \mathbf{x}') = -\eta' \nabla_{\mathbf{x}', \mathbf{z}^*} l_G(\mathbf{x}', \mathbf{z}^*) = \mathbf{B} \end{aligned} \quad (9)$$

Importantly, since $\mathbf{z}_{\text{new}}^0$ is close to \mathbf{z}^* and l_G is smoothed enough (condition (c)), we have

$$\alpha_0 \approx \alpha_1 \dots \approx \alpha_{T'} = \mathbf{A}, \quad \beta_0 \approx \beta_1 \dots \approx \beta_{T'} = \mathbf{B} \quad (10)$$

which yields

$$\nabla_{\mathbf{x}'} \mathbf{z}^{T'} \approx \sum_{t=0}^{T'} \mathbf{B} \mathbf{A}^t \quad (11)$$

substituting $\mathbf{z}^* = \mathbf{z}^{T'}$ completes the proof. \square

Unifying Lemma 1 and Lemma 2 immediately yields Theorem 1. \square

B. Experiments: more details and results

B1. Detailed defense settings

The detailed network architectures for the deep generative models used in DefenseGAN and ABS are summarized in Table 1. Here, C is the number of channels in the image (e.g $C = 3$) and d is the dimensionality of the latent representation.

- $\text{conv}(c, m, k, s, p)$ refers to convolution with c input channels, m feature maps, filter size $k \times k$, stride s and padding p ;
- $\text{conv}^\top(c, m, k, s, p)$ refers to the transpose of convolution (sometimes called deconvolution) (Radford et al., 2015) with c input channels, m feature maps, kernel $k \times k$, stride s and padding p ;

Table 1. Discriminator / Encoder

Stage1	Stage2	Stage3	Stage4	out
conv(C , 256, 4, 2, 1) LeakyReLU(0.2)	conv(256, 512, 4, 2, 1) LeakyReLU(0.2)	conv(512, 1024, 4, 2, 1) LeakyReLU(0.2)	conv(1024, 128, 4, 1, 0) LeakyReLU(0.2)	FC(d)

Table 2. Generator / Decoder

Stage1	Stage2	Stage3	Stage4	out
conv [⊤] (d , 1024, 4, 1, 0) BatchNorm ReLU	conv [⊤] (1024, 512, 4, 2, 1) BatchNorm ReLU	conv [⊤] (512, 256, 4, 2, 1) BatchNorm ReLU	conv [⊤] (256, C , 4, 2, 1)	Tanh

- LeakyReLU(r) refers to the leaky version of ReLU unit (Maas et al., 2013) with parameter r ;
- FC(d) denotes the fully connected layer with d outputs.
- BatchNorm refers to the batch normalization operation (Ioffe & Szegedy, 2015).

All the GANs are realized as Wasserstein GAN (Arjovsky et al., 2017) trained with gradient penalty (Gulrajani et al., 2017).

B2. Detailed attack settings

NES. We estimate the gradient from a population of $n = 50$ samples. These samples are drawn from a search distribution with variance σ^2 set to be $\sigma = 0.5$ in DefenseGAN and $\sigma = 0.1$ in ABS (we discover that the default $\sigma = 0.1$ is not enough to break DefenseGAN in the experiments we consider). We then update the adversarial sample \mathbf{x}' with the estimated gradient by the CW objective:

$$\nabla_{\mathbf{x}'} \mathcal{L}(\mathbf{x}') \approx \nabla_{\mathbf{x}'} \|\mathbf{x}' - \mathbf{x}\|_2^2 + \lambda \cdot \hat{\nabla}_{\mathbf{x}'}^{\text{NES}} F(\mathbf{x}') \quad (12)$$

BPDA. For DefenseGAN (Samangouei et al., 2018), we adopt the same set up as in the original obfuscated gradient literature (Athalye et al., 2018). A CW-like objective is applied so as to achieve minimal distortion. More specifically, it updates the adversarial sample \mathbf{x}' in two steps: (a) forward step, where we find $\mathbf{z}^* = \arg \min_{\mathbf{z}} l_G(\mathbf{z}; \mathbf{x}') = \|G(\mathbf{z}) - \mathbf{x}'\|_2$; (b) backward step, where we update $\mathbf{x}' \leftarrow \mathbf{x}' - \eta \cdot \text{sign}(\nabla_{\mathbf{x}'} \mathcal{L}(\mathbf{x}'))$ with

$$\nabla_{\mathbf{x}'} \mathcal{L}(\mathbf{x}') \approx \nabla_{\mathbf{x}'} \|\mathbf{x}' - \mathbf{x}\|_2^2 + \lambda \cdot \nabla_{\mathbf{x}''} F(\mathbf{x}'') \Big|_{\mathbf{x}' = G(\mathbf{z}^*)} \quad (13)$$

the value of λ are determined by a 5-steps binary search. The value of η is set to be $\eta = 0.04 \cdot t/T$ where t is the current iteration and T is the number of maximum iteration.

For ABS (Schott et al., 2019), we adopt the same ‘latent space attack’ as in the ABS literature (Schott et al., 2019), which is variant of BPDA adapted to this defense. It updates the adversarial sample \mathbf{x}' in two steps: (a) forward step, where we find $\mathbf{z}^* = \arg \min_{\mathbf{z}} l_{G_c}(\mathbf{z}; \mathbf{x}') = \|G_c(\mathbf{z}) - \mathbf{x}'\|_2$ with $c \neq y$ the wrong class; (b) backward step, where we update \mathbf{x}' by

$$\mathbf{x}' \leftarrow (1 - \eta) \cdot \mathbf{x}' + \eta \cdot G_c(\mathbf{z}^*) \quad (14)$$

that is, we make a step towards the best reconstruction in the wrong class. We set $\eta = 0.02$ and iterate until we have found \mathbf{x}' . To ensure that we can find an adversarial sample within the 75 update steps, we increase the value of η every 25 steps. After we have found \mathbf{x}' , a 5-step binary search between \mathbf{x} and \mathbf{x}' is applied to reduce the distortion as much as possible.

B3. Comparison among different defenses

We also compares the robustness of DefenseGAN (Samangouei et al., 2018), ABS (Schott et al., 2019) and adversarial training (ADVT) (Madry et al., 2018). Note that when applied to ADVT our attack degenerates to the original CW attack.

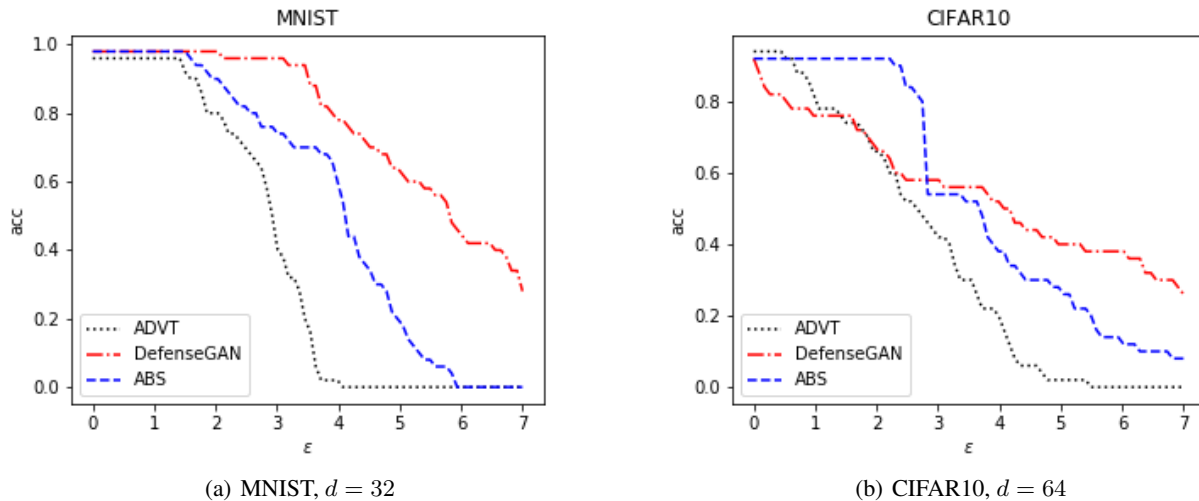


Figure 1. Comparing the robustness of different defenses in the presence of the proposed attack. The horizontal axis denotes the l_2 distortion and the vertical axis denotes the classification accuracy. These results are reported on 100 images only.

Figure 1 shows the accuracy-distortion curve of the various defenses under our attack. Interestingly, we see that despite of our empowered attack, the two deep generative model-based defenses (DefenseGAN, ABS) do seem to be more robust than adversarial training given the experimental setting. A possible explanation is that although DefenseGAN and ABS do not realize the on-manifold conjecture perfectly as analyzed in the paper, they do to some extent exclude some out-off-the-distribution samples, which yield a good protection. We conclude that this kind of defense is still very promising.

References

- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pp. 214–223, 2017.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International Conference on Machine Learning*, pp. 274–283, 2018.
- Gould, S., Fernando, B., Cherian, A., Anderson, P., Cruz, R. S., and Guo, E. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *arXiv preprint arXiv:1607.05447*, 2016.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pp. 5767–5777, 2017.
- Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge university press, 2012.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pp. 448–456, 2015.
- Maas, A. L., Hannun, A. Y., and Ng, A. Y. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, pp. 3, 2013.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- Samangouei, P., Kabkab, M., and Chellappa, R. Defense-gan: Protecting classifiers against adversarial attacks using generative models. In *ICLR*, 2018.

Schott, L., Rauber, J., Bethge, M., and Brendel, W. Towards the first adversarially robust neural network model on mnist. In *ICLR*, pp. 1–16, 2019.

Shaban, A., Cheng, C.-A., Hatch, N., and Boots, B. Truncated back-propagation for bilevel optimization. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1723–1732, 2019.