
Adversarial Attacks on Variational Autoencoders

Anonymous Author(s)

Affiliation

Address

email

1 **Contributions of latex individuals:** all three authors identify as latex.

2 **Contributions of the presenter:** the main idea on how to attack variational autoencoders, how to measure the attack
3 resistance trade-off and its corresponding visualization, partial writing of the final text, editing.

4 Adversarial attacks derail models by crafting malicious inputs. Image classifiers mislabel those inputs — visually
5 indistinguishable from ordinary ones — with high confidence.

6 In comparison to the extensive literature on adversarial attacks for classifiers [1, 2, 3, 4, 5, 6], attacks for autoencoders
7 are mostly unexplored, possibly because those attacks are hard both to perform and to assess [7, 8]. Still, as autoencoders
8 are advanced as powerful schemes for compressing information [9], attacks on them are potentially at least as dangerous
9 as attacks on classifiers.

10 Evaluating generative models is hard [10], there are no clear-cut success criteria for autoencoder reconstruction, and
11 therefore, neither for the attack. We bypass that difficulty by analyzing how inputs and outputs differ across varying
12 compromises between distorting the input and approaching the target. Although, autoencoders admit many variations:
13 sparse [11], denoising [12], variational [13, 14], Wasserstein [15], symmetric [16], etc, we are particularly interested on
14 Variational Autoencoders (VAEs) since they behave as both autoencoders and as generative models, which brought
15 them the community’s attention.

16 Following up on [7], we propose a scheme to attack different VAEs, as well as a quantitative evaluation framework for
17 the attacks that bypass the need for a success criterion. We compare three kinds of autoencoders: simple variational
18 autoencoders (with fully-connected layers), convolutional variational autoencoders, and DRAW — a recently proposed
19 recurrent variational autoencoder with an attention mechanism [17]. We show that the latter is more resistant to the
20 attacks, and that its recurrent and attention mechanism both contribute to the resistance. We run all — statistically
21 validated — experiments in three datasets (MNIST, SVHN, and CelebA) and show that our quantitative assessment
22 correlates well with a qualitative perception of the attacks.

23 Tabacof et al. [7] introduced attacks on autoencoders, showing that they are possible and much harder than attacks on
24 classifiers. They proposed the graphs we call Distortion–Distortion plots here and evaluated attack success by visual
25 inspection of those graphs. Right after, Kos et al. [8] followed up with a work that attacked both the latent representation
26 and the output of VAE–GAN autoencoders.

27 We explore here two types of attacks on VAEs: 1) input attack where the optimization goal is to find the distortion
28 which minimizes the ℓ_2 distance between the target and the VAE’s reconstruction image; and 2) latent attack where the
29 goal is to minimize the Kullback–Leibler divergence between the target’s and the VAE’s resulting latent variables. In
30 both methods, a regularization term is added to the optimization goal in order to keep the distortion norm small.

31 However, there is no sharp criterion to define whether the attack succeeded. We address this shortcoming with the
32 **AUDDC** (Area under Distortion–Distortion Curve). For a given original and target pair, we compute different results,
33 with different approximation compromises. The Distortion–Distortion plots show, for each attempt, how much we
34 distorted the original and how much we approached the target (both measured by ℓ_2). We add limiting lines to the plot:
35 no distortion added (and original reconstruction) at the leftmost/gray and topmost/orange lines; the ℓ_2 -distance between
36 the target and the reconstruction of the target by the model at the bottommost/red line; the ℓ_2 -distance between the
37 original and target image. Those limits represent, respectively, the starting point, the intrinsic limitation of the model,
38 and the maximum “sensible” distortion (which allows going from the original to the target directly). We normalize
39 the graph so that the distance between those lines is 1. The AUDDC is the area under the curve given by the linear
40 interpolation of the points. The closer this area is to 1, the more resistant the model was to the attack (and the less
41 successful the attack was) (Figure 1).

42 We employed three datasets: MNIST [18], SVHN [19], and CelebA [20]. We evaluated four models: VAE with only
43 fully-connected layers (VAE); VAE with (de)convolutional layers (CVAE); the recurrent autoencoder DRAW [17]
44 without and with its attention mechanism and different number of recurrent steps. For each pair model–dataset,
45 we run 20 attacks with different pairs of image–target. For the quantitative analysis, we averaged the AUDDC

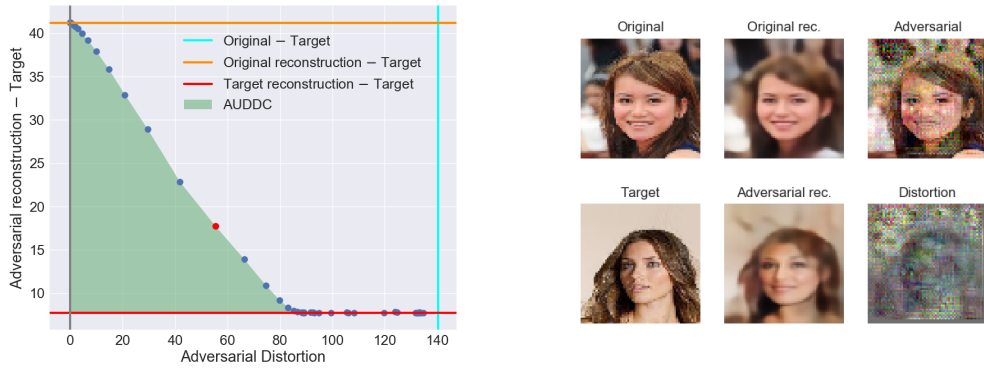


Figure 1: Left: the proposed metric: Area Under the Distortion–Distortion Curve (AUDDC). Right: visualization of a single point (red dot) of the left plot.

Table 1: Average \pm 95%-confidence interval of AUDDC (times 100) for all models and datasets. Higher values indicate higher resistance to the attacks.

Steps	VAE	CVAE	DRAW*	DRAW	DRAW*	DRAW	Average
Attacks on latent representation							
MNIST	27 \pm 2	35 \pm 3	27 \pm 1	35 \pm 3	71 \pm 5	91 \pm 3	47 \pm 3
SVHN	19 \pm 1	18 \pm 1	09 \pm 1	27 \pm 2	74 \pm 6	96 \pm 2	41 \pm 4
CelebA	31 \pm 1	28 \pm 1	21 \pm 2	36 \pm 1	81 \pm 4	97 \pm 1	49 \pm 4
Average	25 \pm 1	27 \pm 2	19 \pm 2	33 \pm 1	75 \pm 3	95 \pm 1	46 \pm 2
Attacks on output							
MNIST	35 \pm 2	56 \pm 3	38 \pm 2	48 \pm 4	29 \pm 3	69 \pm 4	46 \pm 2
SVHN	19 \pm 1	19 \pm 2	13 \pm 1	27 \pm 2	21 \pm 2	34 \pm 2	22 \pm 1
CelebA	27 \pm 1	24 \pm 1	31 \pm 3	35 \pm 1	29 \pm 2	40 \pm 1	31 \pm 1
Average	27 \pm 1	33 \pm 3	27 \pm 2	37 \pm 2	26 \pm 1	47 \pm 3	33 \pm 1
All attacks							
MNIST	31 \pm 2	45 \pm 3	32 \pm 2	42 \pm 3	50 \pm 5	80 \pm 3	47 \pm 2
SVHN	19 \pm 1	19 \pm 1	11 \pm 1	27 \pm 1	47 \pm 7	65 \pm 7	31 \pm 2
CelebA	29 \pm 1	26 \pm 1	26 \pm 2	36 \pm 1	55 \pm 6	68 \pm 7	40 \pm 2
Average	26 \pm 1	30 \pm 2	23 \pm 1	35 \pm 1	51 \pm 4	71 \pm 3	39 \pm 1

* Attention mechanism disabled.

46 for the chosen factors. To check which factors lead to significant influence, we used a multi-way ANOVA, with
 47 second-order interactions, and post-hoc Tukey honest significant differences which found significant differences (all
 48 p-values $<$ 0.015) for all pairs of levels of all factors shown on the Table 1.

49 Attacking auto-encoders is relatively difficult if compared to attacking classifiers, where the distortions can be invisible
 50 to the human eye. Interestingly, DRAW, in particular, was much more resistant to our attacks. No attack succeed in
 51 reconstructing the target image well without incurring in immediately visible distortions to the input. Still, not all
 52 attempts are equal: some models are significantly more resistant than others. The AUDDC metric allows to quantify
 53 that resistance, bypassing the need to establish a clear-cut criterion of success for the attacks, and it correlates well with
 54 the qualitative results.

55 The code to reproduce all experiments will be made available after review.

References

- [1] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014. arXiv:1312.6199.
- [2] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. arXiv:1412.6572.
- [3] Pedro Tabacof and Eduardo Valle. Exploring the space of adversarial images. In *2016 International Joint Conference on Neural Networks, IJCNN 2016, Vancouver, BC, Canada, July 24-29, 2016*, pages 426–433, 2016.
- [4] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *International Conference on Learning Representations (ICLR) Workshop*, 2017. arXiv:1607.02533.
- [5] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method to fool deep neural networks. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2574–2582, 2016.
- [6] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 39–57, 2017.
- [7] Pedro Tabacof, Julia Tavares, and Eduardo Valle. Adversarial images for variational autoencoders. *CoRR*, abs/1612.00155, 2016.
- [8] Jernej Kos, Ian Fischer, and Dawn Song. Adversarial examples for generative models. *CoRR*, abs/1702.06832, 2017.
- [9] Karol Gregor, Frederic Besse, Danilo Jimenez Rezende, Ivo Danihelka, and Daan Wierstra. Towards conceptual compression. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 3549–3557, 2016.
- [10] Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016. arXiv:1511.01844.
- [11] Andrew Ng. Sparse autoencoder. *CS294A Lecture notes*, 72:1–19, 2011.
- [12] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *The Journal of Machine Learning Research*, 11:3371–3408, 2010.
- [13] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014. arXiv:1312.6114.
- [14] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pages 1278–1286, 2014.
- [15] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations*, 2018.
- [16] Yuchen Pu, Weiyao Wang, Ricardo Henao, Liqun Chen, Zhe Gan, Chunyuan Li, and Lawrence Carin. Adversarial symmetric variational autoencoder. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4330–4339. Curran Associates, Inc., 2017.
- [17] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra. DRAW: A recurrent neural network for image generation. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 1462–1471, 2015.
- [18] Yann LeCun, Corinna Cortes, and Christopher JC Burges. The MNIST database of handwritten digits, 1998.
- [19] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS 2011 Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [20] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV), ICCV '15*, pages 3730–3738, Washington, DC, USA, 2015. IEEE Computer Society.