
Likelihood Estimation for Generative Adversarial Networks

Hamid Eghbal-zadeh¹ Gerhard Widmer¹

Abstract

We present a simple method for assessing the quality of generated images in Generative Adversarial Networks (GANs). The method can be applied in any kind of GAN without interfering with the learning procedure or affecting the learning objective. The central idea is to define a likelihood function that correlates with the quality of the generated images. In particular, we derive a Gaussian likelihood function from the distribution of the embeddings (hidden activations) of the real images in the discriminator, and based on this, define two simple measures of how likely it is that the embeddings of generated images are from the distribution of the embeddings of the real images. This yields a simple measure of fitness for generated images, for all varieties of GANs. Empirical results on CIFAR-10 demonstrate a strong correlation between the proposed measures and the perceived quality of the generated images.

1. Introduction

Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) are one of the big discoveries in deep learning in recent years. GANs can learn to generate realistic images through an iterative game between a generator network and a discriminator network. Although GANs are able to learn the Probability Density Function (PDF) underlying a set of real images, they can not explicitly provide a likelihood for (i.e., effectively evaluate) this estimated PDF. In other words, they cannot provide a direct estimate of the probability of an image to come from the true underlying distribution. As an alternative, many different variations of GANs have been proposed where some other measure of goodness is provided to assess the quality of the generated images. For instance, the Wasserstein

GAN (Arjovsky et al., 2017) relies on the *Wasserstein distance* between the activations of the real and the generated images in the discriminator. Other approaches have been developed that not only use a measure of distance in the discriminator, but also consider a *probability measure* in calculating distances. In the Cramer-GAN (Bellemare et al., 2017), the authors use a distance from the family of integral probability metrics (IPM) which measures the distance between probability distributions of generated and real images in the discriminator’s activations. As another example, MMD-GAN (Li et al., 2017) relies on the maximum mean discrepancy (MMD) between two distributions of the real and fake images.

While all these models rely on a measure of distance, they all do so by using their specific distances as an objective for training the GAN. Hence, they can not be ported into other kinds of GANs with different training objectives. In (Salimans et al., 2016), the *inception score* is introduced for comparing the quality of the generated images across different Generative models. This is basically the Inception Model (Szegedy et al., 2016) trained on ImageNet to get both the conditional and marginal label distribution for the generated image. Inception score requires lots of samples (as suggested, $50k$) and also requires a model trained on ImageNet.

In this paper, we propose *Likelihood Estimation for Generative Adversarial Networks (LeGAN)*, a general measure of goodness, based on a likelihood function defined for the discriminator, which can be applied to any kind of GAN. LeGAN does not interfere with the training of GANs and does not need any pre-training. We will attempt to show the correlation between this measure and perceived image quality in an experiment with CIFAR-10 data.

2. Generative Adversarial Networks

In GANs, a discriminator D tries to distinguish between real and generated (fake) images while competing with a generator G that tries to fool the discriminator D by generating realistic images. The loss of the discriminator D in the vanilla GAN (Goodfellow et al., 2014) is defined as:

$$\mathcal{L}_D = \mathbb{E}_{\mathbf{x} \sim P_{real}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim P_{fake}} [\log(1 - D(G(\mathbf{z})))] \quad (1)$$

¹Department of Computational Perception, Johannes Kepler University of Linz, Austria. Correspondence to: Hamid Eghbal-zadeh <hamid.eghbal-zadeh@jku.at>.

(where it is assumed that for a given input image (vector) \mathbf{x} , the discriminator D outputs a number – the estimated probability of the image coming from the real set); and the loss of the generator G is defined as:

$$\mathcal{L}_G = \mathbb{E}_{\mathbf{z} \sim P_{fake}} [-\log D(G(\mathbf{z}))] \quad (2)$$

where \mathbf{z} is a random observation from a distribution Z ; generator G creates a fake image using this \mathbf{z} .

A modified version of the losses defined above using a least squares criterion is introduced in Least Squares GAN (LSGAN) (Mao et al., 2016) as follows:

$$\mathcal{L}_D^{ls} = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim P_{real}} [(D(\mathbf{x}) - 1)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim P_{fake}} [D(G(\mathbf{z}))^2] \quad (3)$$

and the loss of the generator G in LSGAN is defined as:

$$\mathcal{L}_G^{ls} = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim P_{fake}} [(D(G(\mathbf{z})) - 1)^2] \quad (4)$$

None of the losses defined above provide or use a measure that correlates with the quality of the generated images while the training procedure of GANs. In contrast, in the Wasserstein GAN the generator attempts to minimize the Wasserstein distance between $D(x)$ and $D(G(z))$. As this distance approaches zero, the two distributions become more and more similar as the quality of the generated images improves.

In the following section we explain LeGAN, a tool that can be used with any GAN to provide a meaningful measurement that correlates with the quality of the generated images.

3. Likelihood Estimation for Generative Adversarial Networks

In this section, we describe the proposed measurement to assess the quality of generated images in GANs. The underlying idea is very simple: by estimating the distribution of the embeddings of the real images in the discriminator, we define a likelihood function. This function can then measure the likelihood of any arbitrary image based on the image’s discriminator embedding. For the real images this likelihood is expected to be higher than for fake images. When the generated images are no longer distinguishable from the real ones, the likelihood for the generated images should also become similar to the likelihood of real images. Using central limit theorem (Gnedenko et al., 1954), we choose to model the embeddings of the discriminator using a Gaussian distribution with the assumption that our discriminator embeddings will be roughly normally distributed (which actually turns out to be the case in the experiments mentioned in the next section). But this can be replaced with any other distribution.

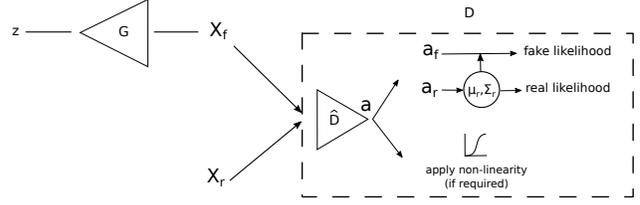


Figure 1. Block diagram of LeGAN. G : the generator. \hat{D} : the discriminator without the final nonlinearity, creating embeddings. D : the discriminator. X_r : real images. X_f : fake (generated) images in one batch of data. μ_r and σ_r^2 : mean and variance of the embeddings of real images. a_r : embeddings of real images. a_f : embeddings of fake images. z : random vector from a known distribution.

A block diagram of LeGAN is provided in Figure 1. The generator G generates (fake) images X_f , given random vectors z from a known distribution as input. The discriminator D , trained with X_f as negative examples and real images X_r as positive examples, tries to distinguish between real and fake samples. In D , a is the embedding¹ representation of images which then can be represented as *probability* by applying activation functions such as *sigmoid*. Hence, a_r represents the embeddings of real images and a_f represents the embeddings of the fake images created by the discriminator network.

Now we define the following likelihood for a given image embedding $a = \hat{D}(\mathbf{x})$ with an assumption of a being Gaussian distributed:

$$\ell_G(a) = P(a | \mathcal{G}) = \mathcal{N}(a; \mu_r, \sigma_r^2) \quad (5)$$

where \hat{D} is the discriminator D without its final activation function. The Gaussian \mathcal{G} is defined by the parameters μ_r (mean) and σ_r^2 (variance) learned from the embeddings of real images in \hat{D} .

Based on the likelihood function $\ell_G(a)$ introduced above, we propose two measurements for GANs:

$$\ell_{diff} = \ell_G(\bar{a}_r) - \ell_G(\bar{a}_f) \quad (6)$$

as the likelihood difference and

$$\ell_{ratio} = \frac{\min(\ell_G(\bar{a}_r), \ell_G(\bar{a}_f))}{\max(\ell_G(\bar{a}_r), \ell_G(\bar{a}_f))} \quad (7)$$

¹These embeddings are the hidden activations of the last layer in the discriminator network before applying the final activation function. Since we want the distribution of the embeddings and not a probability, applying the sigmoid non-linearity is not required.

as the likelihood ratio, where \bar{a}_r and \bar{a}_f are the average of the embeddings in a batch for real and fake images, respectively.

Both of these measures relate to how well the embeddings of fake images match to the Gaussian likelihood function parameterized by the embeddings of real images. When the generated images are close enough to the real images, the distribution of a_r becomes more similar to a_f and therefore a_f fits better in ℓ_G and achieves similar likelihoods as $\ell_G(a_r)$.

4. Distance correlation between images and discriminator embeddings

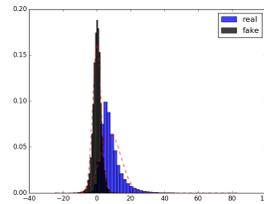
LeGAN provides a measure based on the discriminator’s embedding distribution of real and fake images which correlates with the quality of the generated images. For this to work, we would need to find a way to relate the distance measure in the images to the distance measure between the embeddings of the discriminator. We achieve this property by making our discriminator to be 1-Lipschitz continuous. Assuming x and y two arbitrary images, a discriminator \hat{D} is 1-Lipschitz continuous if:

$$d_{\mathbf{A}}(\hat{D}(x), \hat{D}(y)) \leq d_{\mathbf{X}}(x, y) \quad (8)$$

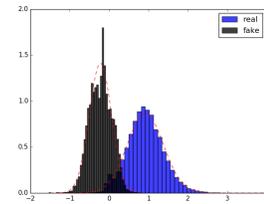
where $d_{\mathbf{A}}$ is a distance defined on \mathbf{A} which is the embedding space of \hat{D} and a distance $d_{\mathbf{X}}$ defined on the images space \mathbf{X} .

To ensure that the embeddings of $\hat{D}(x)$ are bounded, we apply weight-clipping similar to what was proposed in (Arjovsky et al., 2017) to satisfy the 1-Lipschitz continuity of $\hat{D}(x)$ by restricting our weights in the range of $[-0.02, 0.02]$. Since the weights are bounded in this range which is smaller than one, therefore the output of the networks which is sum of the products of the inputs with the weights, is also bounded. Besides, to enforce the network to learn the weights with smaller values, we apply the $L2$ norm penalty on the weights of the discriminator.

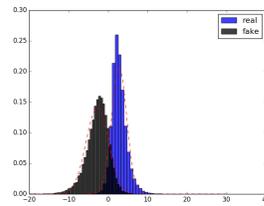
Figure 2 demonstrates the effect of the weight-clipping on the distribution of the discriminator. As can be seen in Figure 2.a and Figure 2.c, the distribution of the embeddings does not become more similar as the model trains and generates better images. In contrast, it can be seen in Figure 2.b and Figure 2.d that at first the real and fake embedding distributions are far apart and as the model trains and generates better images, the distribution of real and fake embeddings becomes more similar and the means of the two clusters become closer. More related examples can be found in the Appendix.



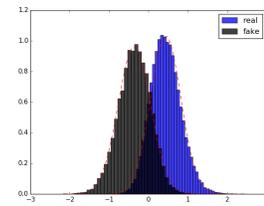
(a) Without weight clipping at epoch 1.



(b) With weight clipping at epoch 1.



(c) Without weight clipping at epoch 13.



(d) With weight clipping at epoch 13.



(e) Generated images at epoch 1.

(f) Generated images at epoch 13.

Figure 2. First and second row: Histograms of real and fake embeddings in $\hat{D}(x)$ of the vanilla GAN (Goodfellow et al., 2014). The third row shows samples of generated images at epoch one and also at epoch 13. Please note the very different ranges in the x and y axis. Third row: generated images.

5. Empirical results

In this section, you can find the empirical results on the similarity measures of LeGAN. Figure 3 shows how the quality of generated images is correlated with the LeGAN measures for an vanilla GAN. Also in Figure 4 another example is provided for an LSGAN to demonstrate this correlation. As can be seen in both example, while the quality of generated images increases with more training updates, ℓ_{diff} decreases which shows that the difference in the likelihood of the real and generated images decreases. By looking at ℓ_{ratio} we can also observe that the ratio between the likelihood of real and fake images increases. For a better understanding of the behavior of LeGAN, we provided more experimental results in the appendix.

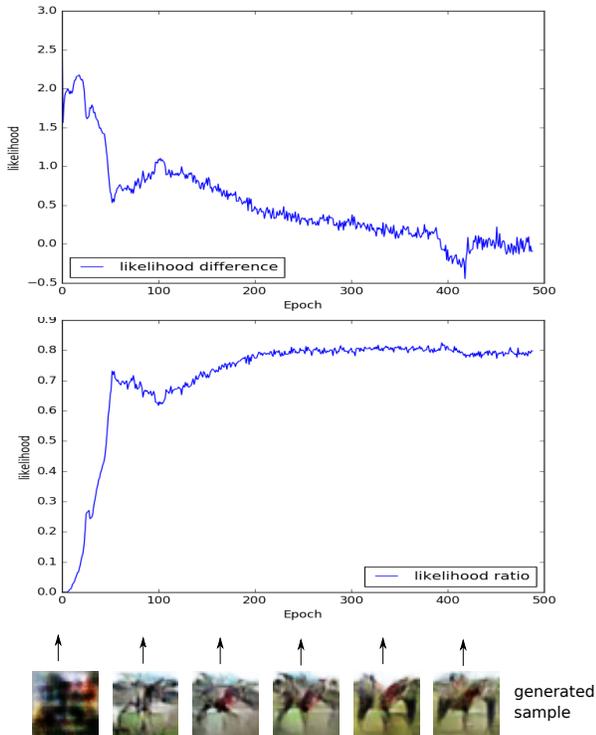


Figure 3. The proposed measures with generated samples using a vanilla GAN(Goodfellow et al., 2014).

6. Conclusion

In this paper, we proposed two measures for the fitness of the generated images based on a Gaussian likelihood function from the distribution of the embeddings of the real images from the discriminator. We provided experimental results CIFAR-10 demonstrating that our measure correlates with the quality of the images generated with any GANs. We used our measurement with the vanilla GAN and LS-

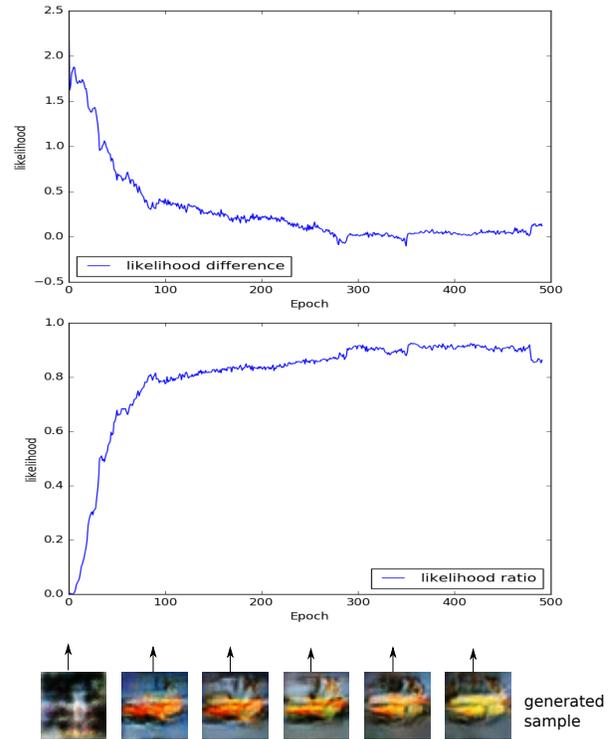


Figure 4. The proposed measures with generated samples using an LSGAN (Mao et al., 2016).

GAN which by default do not have a measure of fitness. We also compared our measures with the Wasserstein distance in Wasserstein GAN², showing our measures show correlation with the Wasserstein distance.

Acknowledgments

This work was supported by the Austrian Ministries BMVIT and BMWFW, and the Province of Upper Austria, via the COMET Center SCCH. We also gratefully acknowledge the support of NVIDIA Corporation with the donation of a Titan X GPU used for this research.

References

- Arjovsky, Martin, Chintala, Soumith, and Bottou, Léon. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- Bastien, Frédéric, Lamblin, Pascal, Pascanu, Razvan, Bergstra, James, Goodfellow, Ian, Bergeron, Arnaud, Bouchard, Nicolas, Warde-Farley, David, and Bengio, Yoshua. Theano: new features and speed improvements.

²These results are provided in the appendix.

arXiv preprint arXiv:1211.5590, 2012.

Bellemare, Marc G, Danihelka, Ivo, Dabney, Will, Mohamed, Shakir, Lakshminarayanan, Balaji, Hoyer, Stephan, and Munos, Rémi. The cramer distance as a solution to biased wasserstein gradients. *arXiv preprint arXiv:1705.10743*, 2017.

Dieleman, Sander, Schltter, Jan, Raffel, Colin, Olson, Eben, Snderby, Sren Kaae, Nouri, Daniel, et al. Lasagne: First release., August 2015. URL <http://dx.doi.org/10.5281/zenodo.27878>.

Gnedenko, BV, Kolmogorov, AN, Gnedenko, BV, and Kolmogorov, AN. Limit distributions for sums of independent. *Am. J. Math.*, 105, 1954.

Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.

Ioffe, Sergey and Szegedy, Christian. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pp. 448–456, 2015.

Kingma, Diederik and Ba, Jimmy. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Li, Chun-Liang, Chang, Wei-Cheng, Cheng, Yu, Yang, Yiming, and Póczos, Barnabás. Mmd gan: Towards deeper understanding of moment matching network. *arXiv preprint arXiv:1705.08584*, 2017.

Maas, Andrew L, Hannun, Awni Y, and Ng, Andrew Y. Rectifier nonlinearities improve neural network acoustic models. In *Proc. ICML*, volume 30, 2013.

Mao, Xudong, Li, Qing, Xie, Haoran, Lau, Raymond YK, Wang, Zhen, and Smolley, Stephen Paul. Least squares generative adversarial networks. *arXiv preprint ArXiv:1611.04076*, 2016.

Salimans, Tim, Goodfellow, Ian, Zaremba, Wojciech, Cheung, Vicki, Radford, Alec, and Chen, Xi. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pp. 2234–2242, 2016.

Szegedy, Christian, Vanhoucke, Vincent, Ioffe, Sergey, Shlens, Jon, and Wojna, Zbigniew. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2818–2826, 2016.

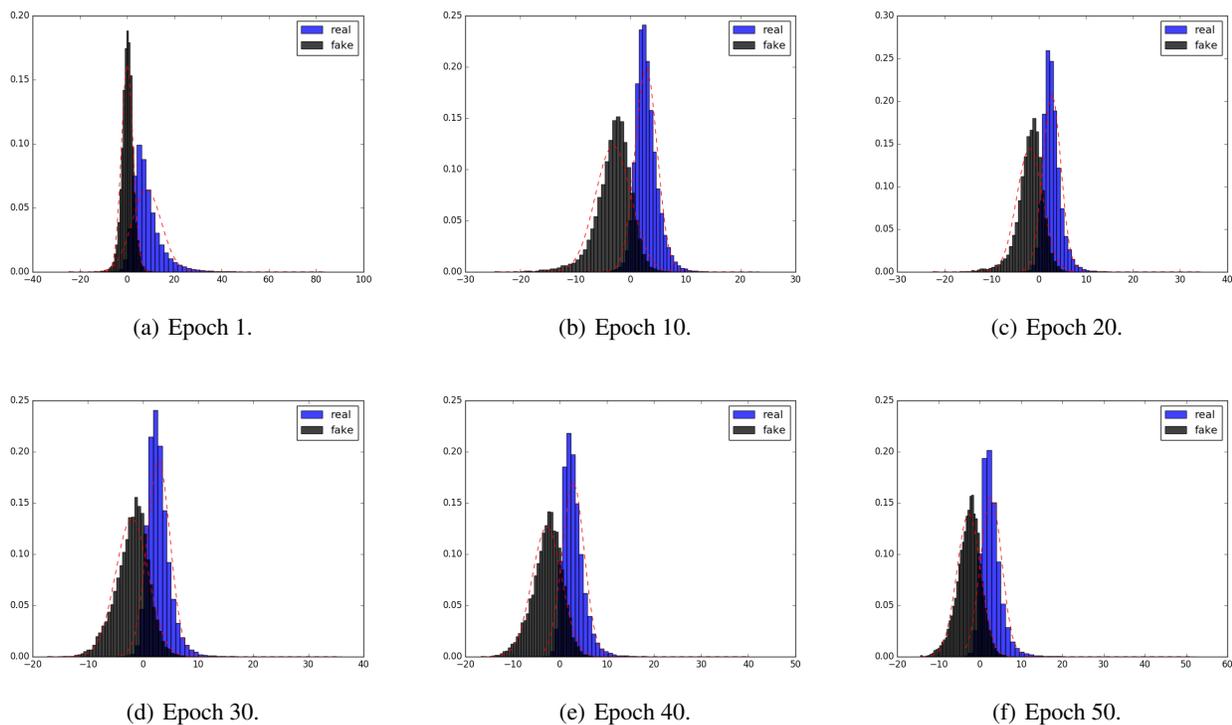


Figure 5. Histograms of the embeddings in different epochs for the vanilla GAN (Goodfellow et al., 2014) without weight-clipping.

7. Appendix

7.1. Extended Empirical Results

In Figure 9, we compare the LeGAN measures with the Wasserstein distance in a Wasserstein GAN. As can be seen, our measures also correlate with the Wasserstein distance as the quality of the generated images improve.

In Figure 5, we provide more examples from the histogram of the discriminator’s embeddings for real and fake images without using the weight-clipping.

Figure 6 shows more examples from the histogram of the discriminator’s embeddings when weight-clipping is applied as discussed in Section 4.

More examples of LeGAN measures for a vanilla GAN can be found in Figure 10. Also in Figure 8, we provide more examples of LeGAN measures with samples of generated images for LSGAN (Mao et al., 2016).

7.2. Architectures

Our neural networks are implemented in Python using the *Lasagne library* (Dieleman et al., 2015) which is based on *Theano* (Bastien et al., 2012).

We observed that applying the weight-clipping reduced the power of our discriminator, therefore we updated our discriminator more often (5 times more than the generator in each epoch).

We used the batch size of 128 and the learning rate of 0.0001 with *ADAM* (Kingma & Ba, 2014) optimizer. We also used RMSProp and achieved similar results. Therefore we only report the results using ADAM.

In Table 1, you can find the architecture of the generator and Table 2 provides the architecture of the discriminator.

Conv: Convolutional layer. BN: Batch-normalization layer (Ioffe & Szegedy, 2015). LReLU: Leaky rectified activation function (Maas et al., 2013).

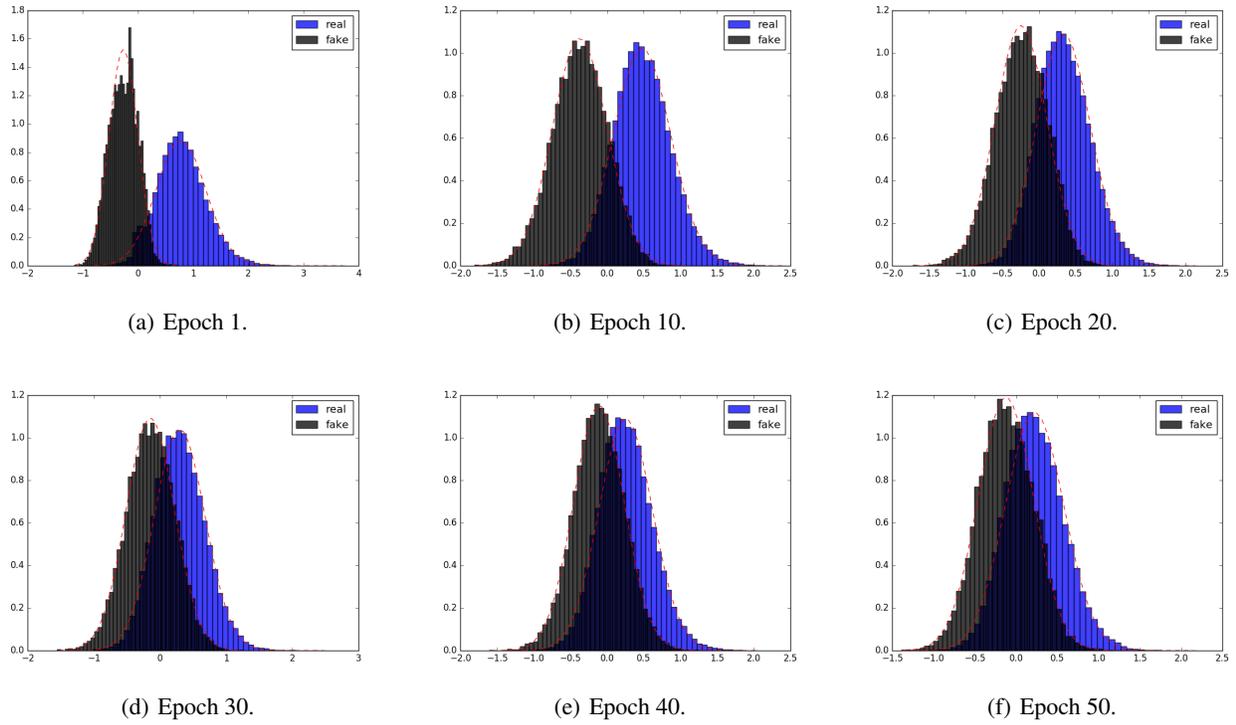


Figure 6. Histograms of the embeddings in different epochs for the vanilla GAN (Goodfellow et al., 2014) with weight-clipping.

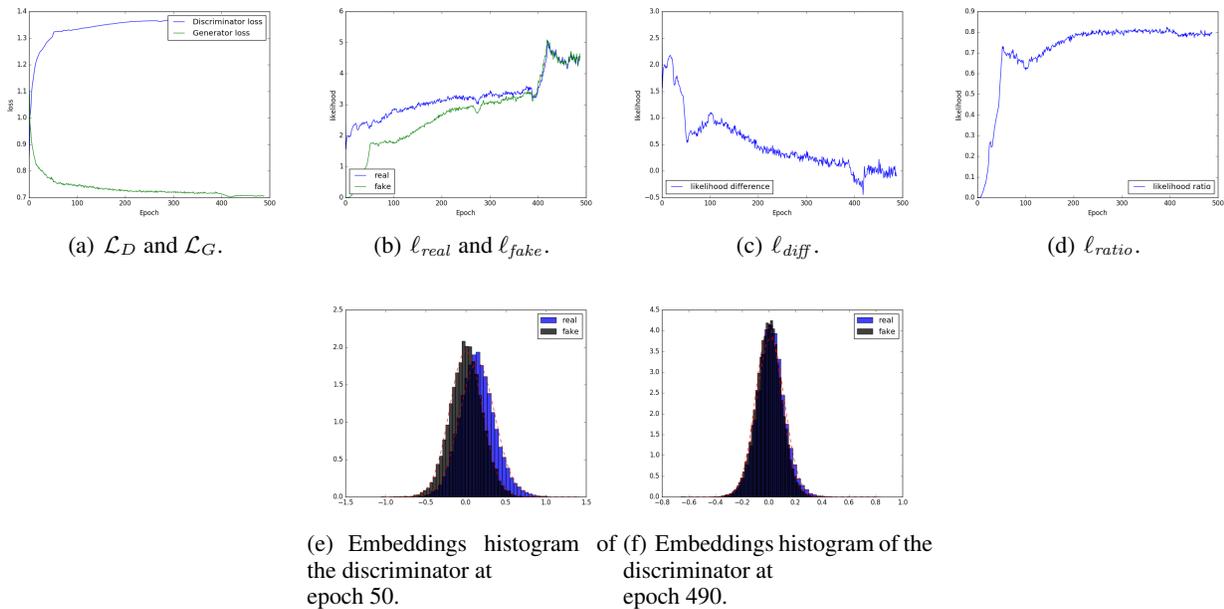


Figure 7. First row: Losses and likelihoods for Discriminator and Generator, LeGAN measures for the vanilla GAN (Goodfellow et al., 2014) with weight-clipping. Second row: histogram of the discriminator embedding.

Likelihood Estimation for Generative Adversarial Networks

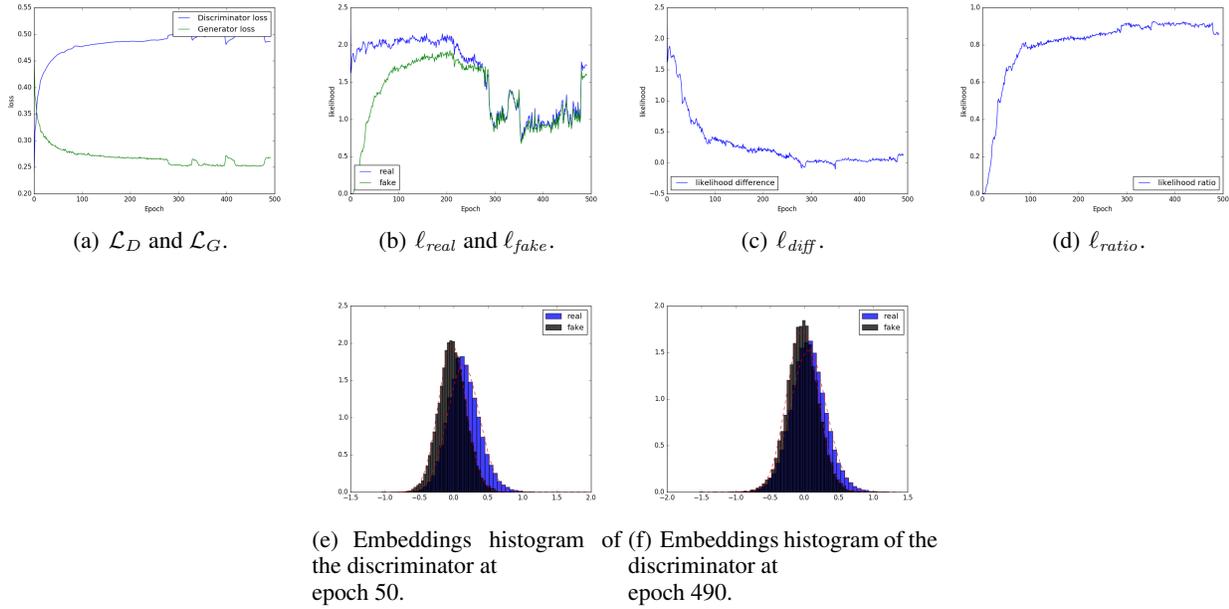


Figure 8. First row: Losses and likelihoods for Discriminator and Generator, LeGAN measures for the LSGAN (Mao et al., 2016) with weight-clipping. Second row: histogram of the discriminator embedding.

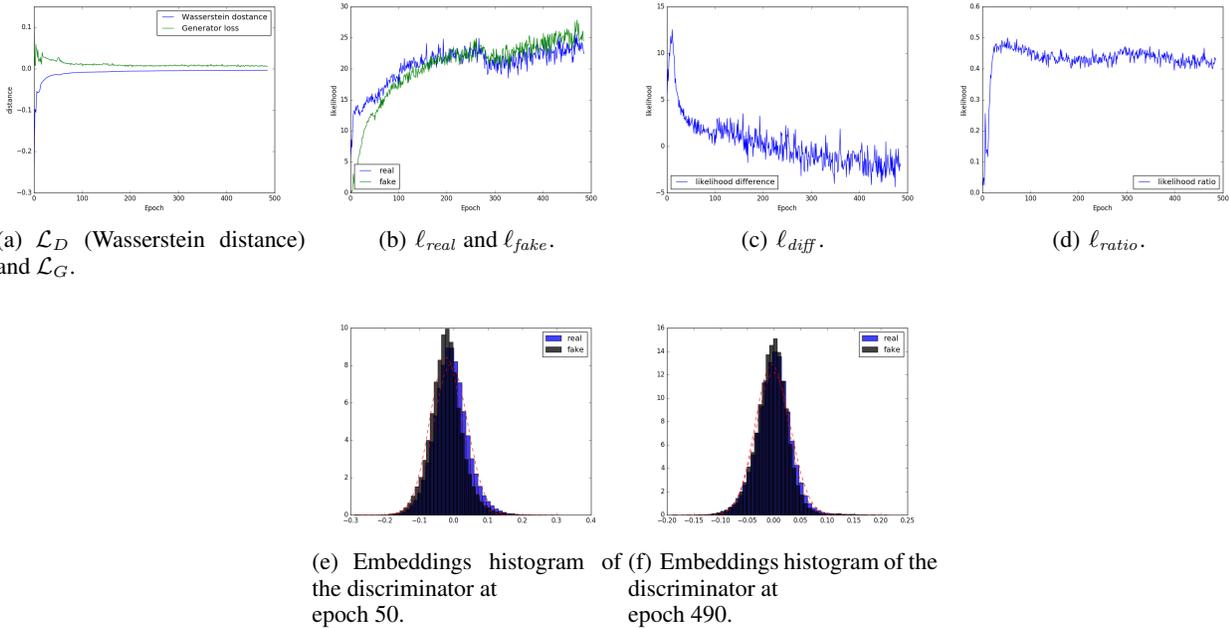


Figure 9. First row: Losses and likelihoods for Discriminator and Generator, LeGAN measures for the Wasserstein GAN (Arjovsky et al., 2017) with weight-clipping. Second row: histogram of the discriminator embedding.



Figure 10. More examples of generated images.

Table 1. The architecture of the Generator.

INPUT 1×100

TRANS CONV LAYER (256, 4×4 , STRIDE=1, CROP=0, NONLIN=LReLU(0.2), INIT=NORMAL(0.02, 0))+BN
 TRANS CONV LAYER (128, 4×4 , STRIDE=2, CROP=1, NONLIN=LReLU(0.2), INIT=NORMAL(0.02, 0))+BN
 TRANS CONV LAYER (64, 4×4 , STRIDE=2, CROP=1, NONLIN=LReLU(0.2), INIT=NORMAL(0.02, 0))+BN
 TRANS CONV LAYER (3, 4×4 , STRIDE=2, CROP=1, NONLIN=SIGMOID, INIT=NORMAL(0.02, 0))+BN

Table 2. The architecture of the Discriminator.

INPUT $3 \times 32 \times 32$

CONV LAYER (64, 4×4 , STRIDE=2, PAD=1, NONLIN=LReLU(0.2), INIT=NORMAL(0.02, 0))
 CONV LAYER (128, 4×4 , STRIDE=2, PAD=1, LReLU(0.2), INIT=NORMAL(0.02, 0))+BN
 CONV LAYER (256, 2×2 , STRIDE=2, PAD=1, LReLU(0.2), INIT=NORMAL(0.02, 0))+BN
 CONV LAYER (1, 4×4 , STRIDE=2, PAD=1, LReLU(0.2)³, INIT=NORMAL(0.02, 0))
