

NIPS 2016 Tutorial: Generative Adversarial Networks

Ian Goodfellow

OpenAI, ian@openai.com

Abstract

This report summarizes the tutorial presented by the author at NIPS 2016 on *generative adversarial networks* (GANs). The tutorial describes: (1) Why generative modeling is a topic worth studying, (2) how generative models work, and how GANs compare to other generative models, (3) the details of how GANs work, (4) research frontiers in GANs, and (5) state-of-the-art image models that combine GANs with other methods. Finally, the tutorial contains three exercises for readers to complete, and the solutions to these exercises.

Introduction

This report¹ summarizes the content of the NIPS 2016 tutorial on *generative adversarial networks* (GANs) (Goodfellow *et al.*, 2014b). The tutorial was designed primarily to ensure that it answered most of the questions asked by audience members ahead of time, in order to make sure that the tutorial would be as useful as possible to the audience. This tutorial is not intended to be a comprehensive review of the field of GANs; many excellent papers are not described here, simply because they were not relevant to answering the most frequent questions, and because the tutorial was delivered as a two hour oral presentation and did not have unlimited time cover all subjects.

The tutorial describes: (1) Why generative modeling is a topic worth studying, (2) how generative models work, and how GANs compare to other generative models, (3) the details of how GANs work, (4) research frontiers in GANs, and (5) state-of-the-art image models that combine GANs with other methods. Finally, the tutorial contains three exercises for readers to complete, and the solutions to these exercises.

The slides for the tutorial are available in PDF and Keynote format at the following URLs:

<http://www.iangoodfellow.com/slides/2016-12-04-NIPS.pdf>

¹This is the arxiv.org version of this tutorial. Some graphics have been compressed to respect arxiv.org's 10MB limit on paper size, and do not reflect the full image quality.

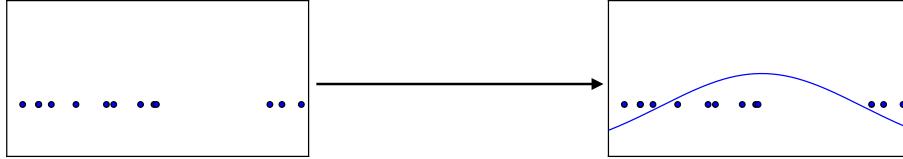


Figure 1: Some generative models perform density estimation. These models take a training set of examples drawn from an unknown data-generating distribution p_{data} and return an estimate of that distribution. The estimate p_{model} can be evaluated for a particular value of \mathbf{x} to obtain an estimate $p_{\text{model}}(\mathbf{x})$ of the true density $p_{\text{model}}(\mathbf{x})$. This figure illustrates the process for a collection of samples of one-dimensional data and a Gaussian model.

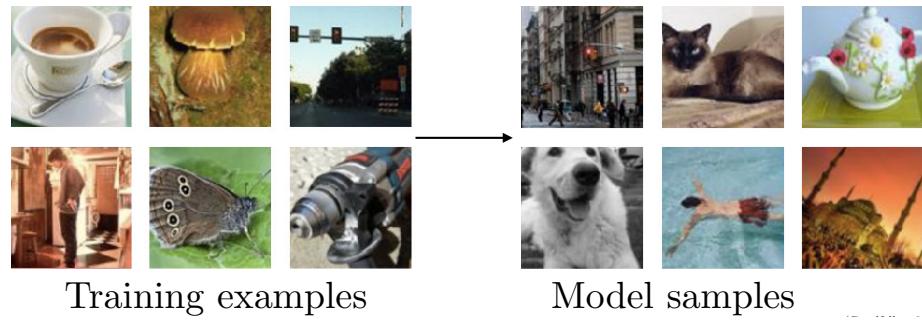


Figure 2: Some generative models are able to generate samples from the model distribution. In this illustration of the process, we show samples from the ImageNet (Deng *et al.*, 2009, 2010; Russakovsky *et al.*, 2014) dataset. An ideal generative model would be able to train on examples as shown on the left and then create more examples from the same distribution as shown on the right. At present, generative models are not yet advanced enough to do this correctly for ImageNet, so for demonstration purposes this figure uses actual ImageNet data to illustrate what an ideal generative model would produce.

<http://www.iangoodfellow.com/slides/2016-12-04-NIPS.key>

The video was recorded by the NIPS foundation and should be made available at a later date.

Generative adversarial networks are an example of *generative models*. The term “generative model” is used in many different ways. In this tutorial, the term refers to any model that takes a training set, consisting of samples drawn from a distribution p_{data} , and learns to represent an estimate of that distribution somehow. The result is a probability distribution p_{model} . In some cases, the model estimates p_{model} explicitly, as shown in figure 1. In other cases, the model is only able to generate samples from p_{model} , as shown in figure 2. Some models are able to do both. GANs focus primarily on sample generation, though it is possible to design GANs that can do both.

1 Why study generative modeling?

One might legitimately wonder why generative models are worth studying, especially generative models that are only capable of generating data rather than providing an estimate of the density function. After all, when applied to images, such models seem to merely provide more images, and the world has no shortage of images.

There are several reasons to study generative models, including:

- Training and sampling from generative models is an excellent test of our ability to represent and manipulate high-dimensional probability distributions. High-dimensional probability distributions are important objects in a wide variety of applied math and engineering domains.
- Generative models can be incorporated into reinforcement learning in several ways. Reinforcement learning algorithms can be divided into two categories; model-based and model-free, with model-based algorithms being those that contain a generative model. Generative models of time-series data can be used to simulate possible futures. Such models could be used for planning and for reinforcement learning in a variety of ways. A generative model used for planning can learn a conditional distribution over future states of the world, given the current state of the world and hypothetical actions an agent might take as input. The agent can query the model with different potential actions and choose actions that the model predicts are likely to yield a desired state of the world. For a recent example of such a model, see Finn *et al.* (2016b), and for a recent example of the use of such a model for planning, see Finn and Levine (2016). Another way that generative models might be used for reinforcement learning is to enable learning in an imaginary environment, where mistaken actions do not cause real damage to the agent. Generative models can also be used to guide exploration by keeping track of how often different states have been visited or different actions have been attempted previously. Generative models, and especially GANs, can also be used for inverse reinforcement learning. Some of these connections to reinforcement learning are described further in section 5.6.
- Generative models can be trained with missing data and can provide predictions on inputs that are missing data. One particularly interesting case of missing data is *semi-supervised learning*, in which the labels for many or even most training examples are missing. Modern deep learning algorithms typically require extremely many labeled examples to be able to generalize well. Semi-supervised learning is one strategy for reducing the number of labels. The learning algorithm can improve its generalization by studying a large number of unlabeled examples which, which are usually easier to obtain. Generative models, and GANs in particular, are able to perform semi-supervised learning reasonably well. This is described further in section 5.4.

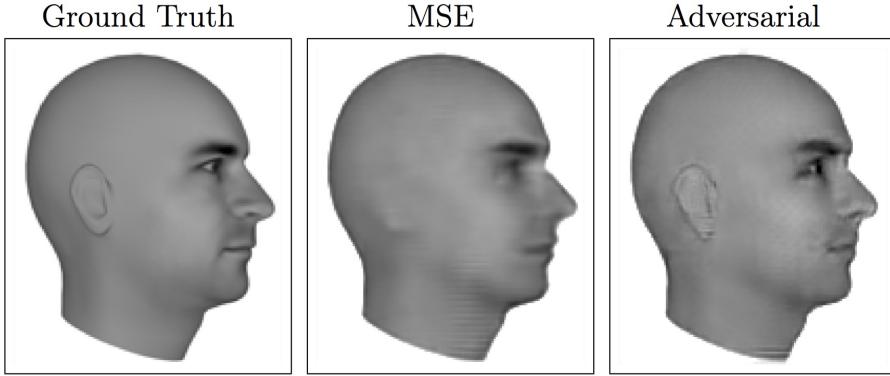


Figure 3: Lotter *et al.* (2015) provide an excellent illustration of the importance of being able to model multi-modal data. In this example, a model is trained to predict the next frame in a video sequence. The video depicts a computer rendering of a moving 3D model of a person’s head. The image on the left shows an example of an actual frame of video, which the model would ideally predict. The image in the center shows what happens when the model is trained using mean squared error between the actual next frame and the model’s predicted next frame. The model is forced to choose a single answer for what the next frame will look like. Because there are many possible futures, corresponding to slightly different positions of the head, the single answer that the model chooses corresponds to an average over many slightly different images. This causes the ears to practically vanish and the eyes to become blurry. Using an additional GAN loss, the image on the right is able to understand that there are many possible outputs, each of which is sharp and recognizable as a realistic, detailed image.

- Generative models, and GANs in particular, enable machine learning to work with *multi-modal* outputs. For many tasks, a single input may correspond to many different correct answers, each of which is acceptable. Some traditional means of training machine learning models, such as minimizing the mean squared error between a desired output and the model’s predicted output, are not able to train models that can produce multiple different correct answers. One example of such a scenario is predicting the next frame in a video, as shown in figure 3.
- Finally, many tasks intrinsically require realitic generation of samples from some distribution.

Examples of some of these tasks that intrinsically require the generation of good samples include:

- *Single image super-resolution:* In this task, the goal is to take a low-resolution image and synthesize a high-resolution equivalent. Generative modeling is required because this task requires the model to impute more information into the image than was originally there in the input.

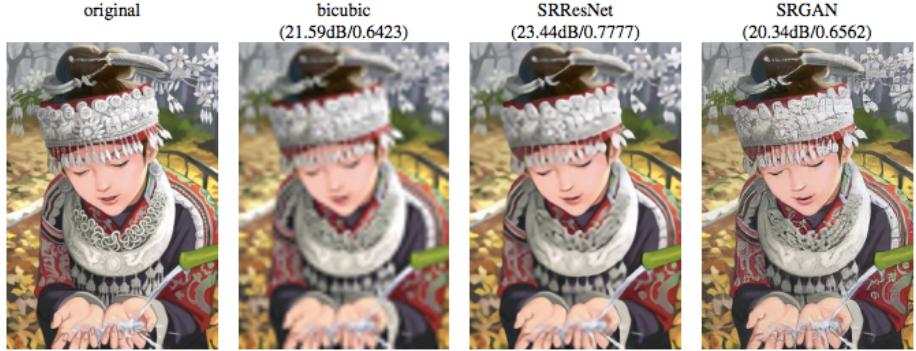


Figure 4: Ledig *et al.* (2016) demonstrate excellent single-image superresolution results that show the benefit of using a generative model trained to generate realistic samples from a multimodal distribution. The leftmost image is an original high-resolution image. It is then downsampled to make a low-resolution image, and different methods are used to attempt to recover the high-resolution image. The bicubic method is simply an interpolation method that does not use the statistics of the training set at all. SRResNet is a neural network trained with mean squared error. SRGAN is a GAN-based neural network that improves over SRGAN because it is able to understand that there are multiple correct answers, rather than averaging over many answers to impose a single best output.

There are many possible high-resolution images corresponding to the low-resolution image. The model should choose an image that is a sample from the probability distribution over possible images. Choosing an image that is the average of all possible images would yield a result that is too blurry to be pleasing. See figure 4.

- Tasks where the goal is to create art. Two recent projects have both demonstrated that generative models, and in particular, GANs, can be used to create interactive programs that assist the user in creating realistic images that correspond to rough scenes in the user’s imagination. See figure 5 and figure 6.
- Image-to-image translation applications can convert aerial photos into maps or convert sketches to images. There is a very long tail of creative applications that are difficult to anticipate but useful once they have been discovered. See figure 7.

All of these and other applications of generative models provide compelling reasons to invest time and resources into improving generative models.



Figure 5: Zhu *et al.* (2016) developed an interactive application called *interactive generative adversarial networks* (iGAN). A user can draw a rough sketch of an image, and iGAN uses a GAN to produce the most similar realistic image. In this example, a user has scribbled a few green lines that iGAN has converted into a grassy field, and the user has drawn a black triangle that iGAN has turned into a detailed mountain. Applications that create art are one of many reasons to study generative models that create images. A video demonstration of iGAN is available at the following URL: <https://www.youtube.com/watch?v=9c4z6YsBGQ0>

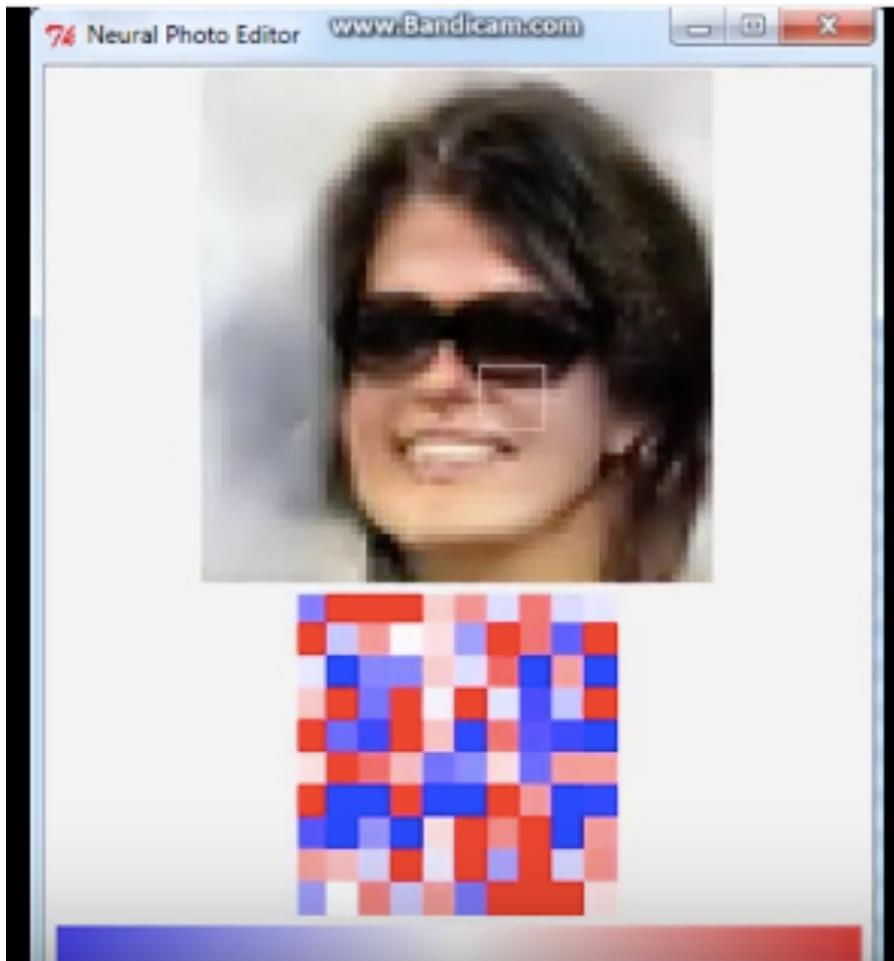


Figure 6: Brock *et al.* (2016) developed *introspective adversarial networks* (IAN). The user paints rough modifications to a photo, such as painting with black paint in an area where the user would like to add black hair, and IAN turns these rough paint strokes into photorealistic imagery matching the user’s desires. Applications that enable a user to make realistic modifications to photo media are one of many reasons to study generative models that create images. A video demonstration of IAN is available at the following URL: <https://www.youtube.com/watch?v=FDELBFSeqQs>



Figure 7: Isola *et al.* (2016) created a concept they called image to image translation, encompassing many kinds of transformations of an image: converting a satellite photo into a map, converting a sketch into a photorealistic image, etc. Because many of these conversion processes have multiple correct outputs for each input, it is necessary to use generative modeling to train the model correctly. In particular, Isola *et al.* (2016) use a GAN. Image to image translation provides many examples of how a creative algorithm designer can find several unanticipated uses for generative models. In the future, presumably many more such creative uses will be found.

2 How do generative models work? How do GANs compare to others?

We now have some idea of what generative models can do and why it might be desirable to build one. Now we can ask: how does a generative model actually work? And in particular, how does a GAN work, in comparison to other generative models?

2.1 Maximum likelihood estimation

To simplify the discussion somewhat, we will focus on generative models that work via the principle of **maximum likelihood**. Not every generative model uses maximum likelihood. Some generative models do not use maximum likelihood by default, but can be made to do so (GANs fall into this category). By ignoring those models that do not use maximum likelihood, and by focusing on the maximum likelihood version of models that do not usually use maximum likelihood, we can eliminate some of the more distracting differences between different models.

The basic idea of maximum likelihood is to define a model that provides an estimate of a probability distribution, parameterized by parameters θ . We then refer to the **likelihood** as the probability that the model assigns to the training data: $\prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta)$, for a dataset containing m training examples $\mathbf{x}^{(i)}$.

The principle of maximum likelihood simply says to choose the parameters for the model that maximize the likelihood of the training data. This is easiest to do in log space, where we have a sum rather than a product over examples. This sum simplifies the algebraic expressions for the derivatives of the likelihood with respect to the models, and when implemented on a digital computer, is less prone to numerical problems, such as underflow resulting from multiplying together several very small probabilities.

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \quad (1)$$

$$= \arg \max_{\boldsymbol{\theta}} \log \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \quad (2)$$

$$= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}). \quad (3)$$

In equation 2, we have used the property that $\arg \max_v f(v) = \arg \max_v \log f(v)$ for positive v , because the logarithm is a function that increases everywhere and does not change the location of the maximum.

The maximum likelihood process is illustrated in figure 8.

We can also think of maximum likelihood estimation as minimizing the **KL divergence** between the data generating distribution and the model:

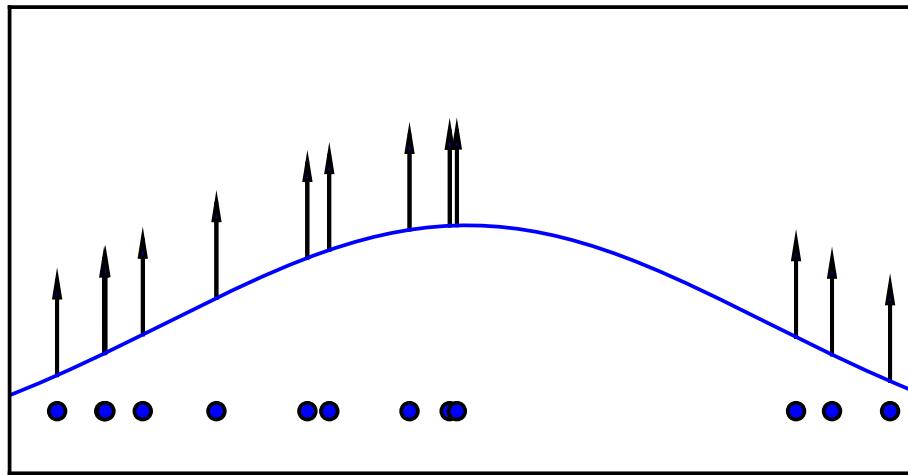
$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})). \quad (4)$$

If we were able to do this precisely, then if p_{data} lies within the family of distributions $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$, the model would recover p_{data} exactly. In practice, we do not have access to p_{data} itself, but only to a training set consisting of m samples from p_{data} . We use these to define \hat{p}_{data} , an **empirical distribution** that places mass only on exactly those m points, approximating p_{data} . Minimizing the KL divergence between \hat{p}_{data} and p_{model} is exactly equivalent to maximizing the log-likelihood of the training set.

For more information on maximum likelihood and other statistical estimators, see chapter 5 of Goodfellow *et al.* (2016).

2.2 A taxonomy of deep generative models

If we restrict our attention to deep generative models that work by maximizing the likelihood, we can compare several models by contrasting the ways that they compute either the likelihood and its gradients, or approximations to these quantities. As mentioned earlier, many of these models are often used with principles other than maximum likelihood, but we can examine the maximum likelihood variant of each of them in order to reduce the amount of distracting differences between the methods. Following this approach, we construct the taxonomy shown in figure 9. Every leaf in this taxonomic tree has some



$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} \log p_{\text{model}}(x \mid \theta)$$

Figure 8: The maximum likelihood process consists of taking several samples from the data generating distribution to form a training set, then pushing up on the probability the model assigns to those points, in order to maximize the likelihood of the training data. This illustration shows how different data points push up on different parts of the density function for a Gaussian model applied to 1-D data. The fact that the density function must sum to 1 means that we cannot simply assign infinite likelihood to all points; as one point pushes up in one place it inevitably pulls down in other places. The resulting density function balances out the upward forces from all the data points in different locations.

advantages and disadvantages. GANs were designed to avoid many of the disadvantages present in pre-existing nodes of the tree, but also introduced some new disadvantages.

2.3 Explicit density models

In the left branch of the taxonomy shown in figure 9 are models that define an explicit density function $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$. For these models, maximization of the likelihood is straightforward; we simply plug the model’s definition of the density function into the expression for the likelihood, and follow the gradient uphill.

The main difficulty present in explicit density models is designing a model that can capture all of the complexity of the data to be generated while still maintaining computational tractability. There are two different strategies used to confront this challenge: (1) careful construction of models whose structure guarantees their tractability, as described in section 2.3.1, and (2) models that admit tractable approximations to the likelihood and its gradients, as described in section 2.3.2.

2.3.1 Tractable explicit models

In the leftmost leaf of the taxonomic tree of figure 9 are the models that define an explicit density function that is computationally tractable. There are currently two popular approaches to tractable explicit density models: fully visible belief networks and nonlinear independent components analysis.

Fully visible belief networks Fully visible belief networks (Frey *et al.*, 1996; Frey, 1998) or FVBMs are models that use the chain rule of probability to decompose a probability distribution over an n -dimensional vector \mathbf{x} into a product of one-dimensional probability distributions:

$$p_{\text{model}}(\mathbf{x}) = \prod_{i=1}^n p_{\text{model}}(x_i | x_1, \dots, x_{i-1}). \quad (5)$$

FVBMs are, as of this writing, one of the three most popular approaches to generative modeling, alongside GANs and variational autoencoders. They form the basis for sophisticated generative models from DeepMind, such as WaveNet (Oord *et al.*, 2016). WaveNet is able to generate realistic human speech. The main drawback of FVBMs is that samples must be generated one entry at a time: first x_1 , then x_2 , etc., so the cost of generating a sample is $O(n)$. In modern FVBMs such as WaveNet, the distribution over each x_i is computed by a deep neural network, so each of these n steps involves a nontrivial amount of computation. Moreover, these steps cannot be parallelized. WaveNet thus requires two minutes of computation time to generate one second of audio, and cannot yet be used for interactive conversations. GANs were designed to be able to generate all of \mathbf{x} in parallel, yielding greater generation speed.

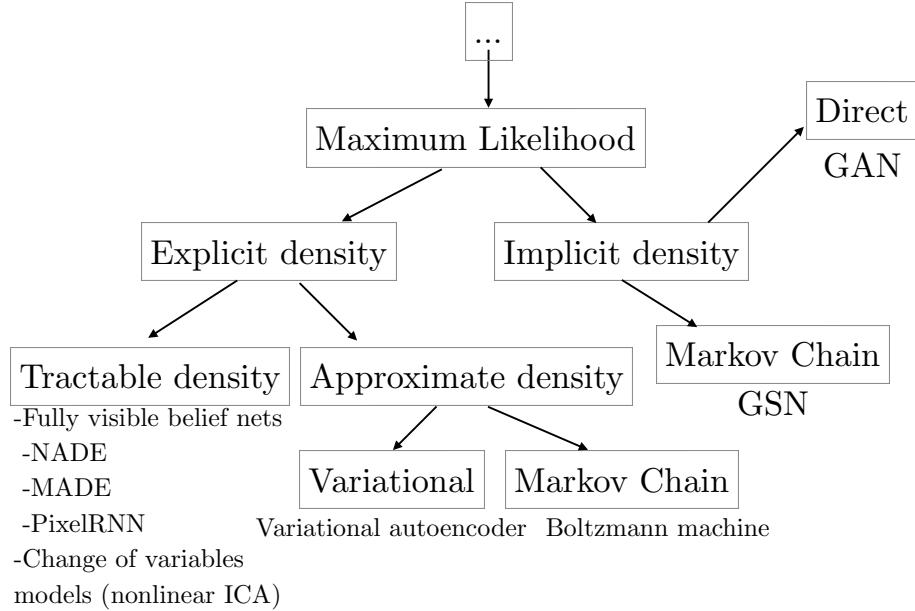


Figure 9: Deep generative models that can learn via the principle of maximum likelihood differ with respect to how they represent or approximate the likelihood. On the left branch of this taxonomic tree, models construct an explicit density, $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$, and thus an explicit likelihood which can be maximized. Among these explicit density models, the density may be computationally tractable, or it may be intractable, meaning that to maximize the likelihood it is necessary to make either **variational approximations** or **Monte Carlo approximations** (or both). On the right branch of the tree, the model does not explicitly represent a probability distribution over the space where the data lies. Instead, the model provides some way of interacting less directly with this probability distribution. Typically the indirect means of interacting with the probability distribution is the ability to draw samples from it. Some of these implicit models that **offer the ability to sample from the distribution do so using a Markov Chain**; the model defines a way to stochastically transform an existing sample in order to obtain another sample from the same distribution. Others are able to generate a sample in a single step, starting without any input. While the models used for GANs can sometimes be constructed to define an explicit density, the training algorithm for GANs makes use only of the model's ability to generate samples. GANs are thus trained using the strategy from the rightmost leaf of the tree: using an implicit model that samples directly from the distribution represented by the model.



Figure 10: Samples generated by a real NVP model trained on 64x64 ImageNet images. Figure reproduced from Dinh *et al.* (2016).

Nonlinear independent components analysis Another family of deep generative models with explicit density functions is based on defining continuous, nonlinear transformations between two different spaces. For example, if there is a vector of latent variables \mathbf{z} and a continuous, differentiable, invertible transformation g such that $g(\mathbf{z})$ yields a sample from the model in \mathbf{x} space, then

$$p_x(\mathbf{x}) = p_z(g^{-1}(\mathbf{x})) \left| \det \left(\frac{\partial g^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|. \quad (6)$$

The density p_x is tractable if the density p_z is tractable and the determinant of the Jacobian of g^{-1} is tractable. In other words, a simple distribution over \mathbf{z} combined with a transformation g that warps space in complicated ways can yield a complicated distribution over \mathbf{x} , and if g is carefully designed, the density is tractable too. Models with nonlinear g functions date back at least to Deco and Brauer (1995). The latest member of this family is real NVP (Dinh *et al.*, 2016). See figure 10 for some visualizations of ImageNet samples generated by real NVP. The main drawback to nonlinear ICA models is that they impose restrictions on the choice of the function g . In particular, the invertibility requirement means that the latent variables \mathbf{z} must have the same dimensionality as \mathbf{x} . GANs were designed to impose very few requirements on g , and, in particular, admit the use of \mathbf{z} with larger dimension than \mathbf{x} .

For more information about the chain rule of probability used to define FVBMs or about the effect of deterministic transformations on probability densities as used to define nonlinear ICA models, see chapter 3 of Goodfellow *et al.* (2016).

In summary, models that define an explicit, tractable density are highly effective, because they permit the use of an optimization algorithm directly on the log-likelihood of the training data. However, the family of models that provide a tractable density is limited, with different families having different disadvantages.

2.3.2 Explicit models requiring approximation

To avoid some of the disadvantages imposed by the design requirements of models with tractable density functions, other models have been developed that still provide an explicit density function but use one that is intractable, requiring the use of approximations to maximize the likelihood. These fall roughly into two categories: those using **deterministic approximations**, which almost always means variational methods, and those using **stochastic approximations**, meaning Markov chain Monte Carlo methods.

Variational approximations Variational methods define a lower bound

$$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}) \leq \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}). \quad (7)$$

A learning algorithm that maximizes \mathcal{L} is guaranteed to obtain at least as high a value of the log-likelihood as it does of \mathcal{L} . For many families of models, it is possible to define an \mathcal{L} that is computationally tractable even when the log-likelihood is not. Currently, the most popular approach to variational learning in deep generative models is the **variational autoencoder** (Kingma, 2013; Rezende *et al.*, 2014) or VAE. Variational autoencoders are one of the three approaches to deep generative modeling that are the most popular as of this writing, along with FVBMs and GANs. The main drawback of variational methods is that, when **too weak of an approximate posterior distribution or too weak of a prior distribution is used**,² even with a perfect optimization algorithm and infinite training data, the gap between \mathcal{L} and the true likelihood can result in p_{model} learning something other than the true p_{data} . GANs were designed to be unbiased, in the sense that with a large enough model and infinite data, the Nash equilibrium for a GAN game corresponds to recovering p_{data} exactly. In practice, variational methods often obtain very good likelihood, but are regarded as producing lower quality samples. There is not a good method of quantitatively measuring sample quality, so this is a subjective opinion, not an empirical fact. See figure 11 for an example of some samples drawn from a VAE. While it is difficult to point to a single aspect of GAN design and say that it results in better sample quality, GANs are generally regarded as producing better samples. Compared to FVBMs, VAEs are regarded as more difficult to optimize, but GANs are not an improvement in this respect. For more information about variational approximations, see chapter 19 of Goodfellow *et al.* (2016).

² Empirically, VAEs with highly flexible priors or highly flexible approximate posteriors can obtain values of \mathcal{L} that are near their own log-likelihood (Kingma *et al.*, 2016; Chen *et al.*, 2016b). Of course, this is testing the gap between the objective and the bound at the maximum of the bound; it would be better, but not feasible, to test the gap at the maximum of the objective. VAEs obtain likelihoods that are competitive with other methods, suggesting that they are also near the maximum of the objective. In personal conversation, L. Dinh and D. Kingma have conjectured that a family of models (Dinh *et al.*, 2014; Rezende and Mohamed, 2015; Kingma *et al.*, 2016; Dinh *et al.*, 2016) usable as VAE priors or approximate posteriors are universal approximators. If this could be proven, it would establish VAEs as being asymptotically consistent.



Figure 11: Samples drawn from a VAE trained on the CIFAR-10 dataset. Figure reproduced from Kingma *et al.* (2016).

Markov chain approximations Most deep learning algorithms make use of some form of stochastic approximation, at the very least in the form of using a small number of randomly selected training examples to form a minibatch used to minimize the expected loss. Usually, sampling-based approximations work reasonably well as long as a fair sample can be generated quickly (e.g. selecting a single example from the training set is a cheap operation) and as long as the variance across samples is not too high. Some models require the generation of more expensive samples, using Markov chain techniques. A Markov chain is a process for generating samples by repeatedly drawing a sample $\mathbf{x}' \sim q(\mathbf{x}' | \mathbf{x})$. By repeatedly updating \mathbf{x} according to the transition operator q , Markov chain methods can sometimes guarantee that \mathbf{x} will eventually converge to a sample from $p_{\text{model}}(\mathbf{x})$. Unfortunately, this convergence can be very slow, and there is no clear way to test whether the chain has converged, so in practice one often uses \mathbf{x} too early, before it has truly converged to be a fair sample from p_{model} . In high-dimensional spaces, Markov chains become less efficient. Boltzmann machines (Fahlman *et al.*, 1983; Ackley *et al.*, 1985; Hinton *et al.*, 1984; Hinton and Sejnowski, 1986) are a family of generative models that rely on Markov chains both to train the model or to generate a sample from the model. Boltzmann machines were an important part of the deep learning renaissance beginning in 2006 (Hinton *et al.*, 2006; Hinton, 2007) but they are now used only very rarely, presumably mostly because the underlying Markov chain approximation techniques have not scaled to problems like ImageNet generation. Moreover, even if Markov chain methods scaled well enough to be used for training, the use of a Markov chain to generate samples from a trained model is undesirable compared to single-step generation methods because the multi-step Markov chain approach has higher computational cost. GANs were designed to avoid using Markov chains for these reasons. For more information about Markov chain Monte Carlo approximations, see chapter 18 of Goodfellow *et al.* (2016). For more information about Boltzmann machines, see chapter 20 of the same book.

Some models use both variational and Markov chain approximations. For example, deep Boltzmann machines make use of both types of approximation (Salakhutdinov and Hinton, 2009).

2.4 Implicit density models

Some models can be trained without even needing to explicitly define a density functions. These models instead offer a way to train the model while interacting only indirectly with p_{model} , usually by sampling from it. These constitute the second branch, on the right side, of our taxonomy of generative models depicted in figure 9.

Some of these implicit models based on drawing samples from p_{model} define a Markov chain transition operator that must be run several times to obtain a sample from the model. From this family, the primary example is the **generative stochastic network** (Bengio *et al.*, 2014). As discussed in section 2.3.2, Markov chains often fail to scale to high dimensional spaces, and impose increased computational costs for using the generative model. GANs

were designed to avoid these problems.

Finally, the rightmost leaf of our taxonomic tree is the family of implicit models that can generate a sample in a single step. At the time of their introduction, GANs were the only notable member of this family, but since then they have been joined by additional models based on kernelized moment matching (Li *et al.*, 2015; Dziugaite *et al.*, 2015).

2.5 Comparing GANs to other generative models

In summary, GANs were designed to avoid many disadvantages associated with other generative models:

- They can generate samples in parallel, instead of using runtime proportional to the dimensionality of \mathbf{x} . This is an advantage relative to FVBMs.
- The design of the generator function has very few restrictions. This is an advantage relative to Boltzmann machines, for which few probability distributions admit tractable Markov chain sampling, and relative to nonlinear ICA, for which the generator must be invertible and the latent code \mathbf{z} must have the same dimension as the samples \mathbf{x} .
- No Markov chains are needed. This is an advantage relative to Boltzmann machines and GSNs.
- No variational bound is needed, and specific model families usable within the GAN framework are already known to be universal approximators, so GANs are already known to be asymptotically consistent. Some VAEs are conjectured to be asymptotically consistent, but this is not yet proven.
- GANs are subjectively regarded as producing better samples than other methods.

At the same time, GANs have taken on a new disadvantage: training them requires finding the Nash equilibrium of a game, which is a more difficult problem than optimizing an objective function.

3 How do GANs work?

We have now seen several other generative models and explained that GANs do not work in the same way that they do. But how do GANs themselves work?

3.1 The GAN framework

The basic idea of GANs is to set up a game between two players. One of them is called the **generator**. The generator creates samples that are intended to come from the same distribution as the training data. The other player is the **discriminator**. The discriminator examines samples to determine whether they

are real or fake. The discriminator learns using traditional supervised learning techniques, dividing inputs into two classes (real or fake). The generator is trained to fool the discriminator. We can think of the generator as being like a counterfeiter, trying to make fake money, and the discriminator as being like police, trying to allow legitimate money and catch counterfeit money. To succeed in this game, the counterfeiter must learn to make money that is indistinguishable from genuine money, and the generator network must learn to create samples that are drawn from the same distribution as the training data. The process is illustrated in figure 12.

Formally, GANs are a structured probabilistic model (see chapter 16 of Goodfellow *et al.* (2016) for an introduction to structured probabilistic models) containing latent variables \mathbf{z} and observed variables \mathbf{x} . The graph structure is shown in figure 13.

The two players in the game are represented by two functions, each of which is differentiable both with respect to its inputs and with respect to its parameters. The discriminator is a function D that takes \mathbf{x} as input and uses $\boldsymbol{\theta}^{(D)}$ as parameters. The generator is defined by a function G that takes \mathbf{z} as input and uses $\boldsymbol{\theta}^{(G)}$ as parameters.

Both players have cost functions that are defined in terms of both players' parameters. The discriminator wishes to minimize $J^{(D)}(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)})$ and must do so while controlling only $\boldsymbol{\theta}^{(D)}$. The generator wishes to minimize $J^{(G)}(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)})$ and must do so while controlling only $\boldsymbol{\theta}^{(G)}$. Because each player's cost depends on the other player's parameters, but each player cannot control the other player's parameters, this scenario is most straightforward to describe as a game rather than as an optimization problem. The solution to an optimization problem is a (local) minimum, a point in parameter space where all neighboring points have greater or equal cost. The solution to a game is a Nash equilibrium. Here, we use the terminology of local differential Nash equilibria (Ratliff *et al.*, 2013). In this context, a Nash equilibrium is a tuple $(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)})$ that is a local minimum of $J^{(D)}$ with respect to $\boldsymbol{\theta}^{(D)}$ and a local minimum of $J^{(G)}$ with respect to $\boldsymbol{\theta}^{(G)}$.

The generator The generator is simply a differentiable function G . When \mathbf{z} is sampled from some simple prior distribution, $G(\mathbf{z})$ yields a sample of \mathbf{x} drawn from p_{model} . Typically, a deep neural network is used to represent G . Note that the inputs to the function G do not need to correspond to inputs to the first layer of the deep neural net; inputs may be provided at any point throughout the network. For example, we can partition \mathbf{z} into two vectors $\mathbf{z}^{(1)}$ and $\mathbf{z}^{(2)}$, then feed $\mathbf{z}^{(1)}$ as input to the first layer of the neural net and add $\mathbf{z}^{(2)}$ to the last layer of the neural net. If $\mathbf{z}^{(2)}$ is Gaussian, this makes \mathbf{x} conditionally Gaussian given $\mathbf{z}^{(1)}$. Another popular strategy is to apply additive or multiplicative noise to hidden layers or concatenate noise to hidden layers of the neural net. Overall, we see that there are very few restrictions on the design of the generator net. If we want p_{model} to have full support on \mathbf{x} space we need the dimension of \mathbf{z} to be at least as large as the dimension of \mathbf{x} , and G

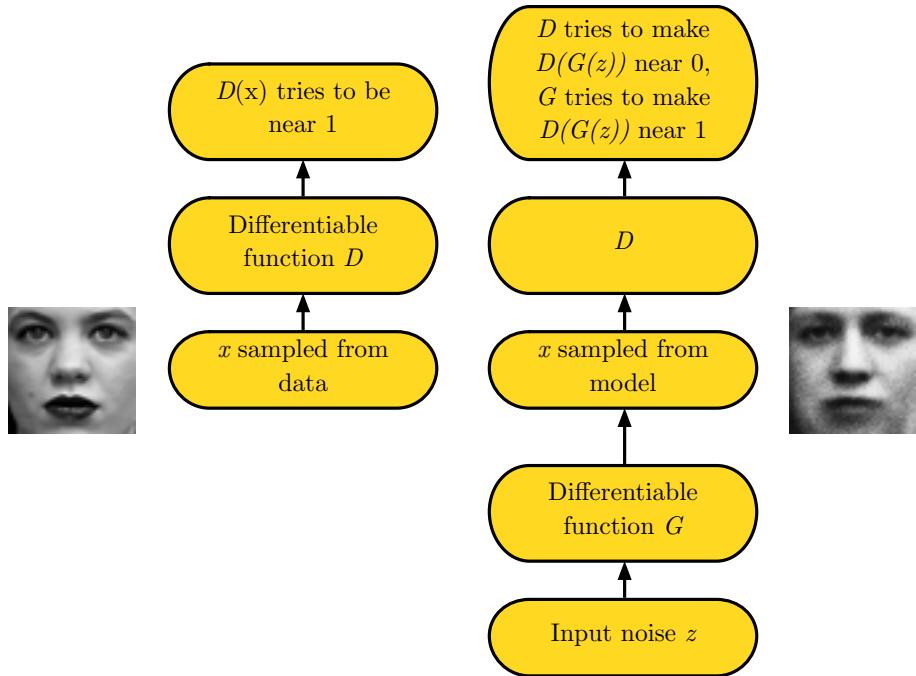


Figure 12: The GAN framework pits two adversaries against each other in a game. Each player is represented by a differentiable function controlled by a set of parameters. Typically these functions are implemented as deep neural networks. The game plays out in two scenarios. In one scenario, training examples \mathbf{x} are randomly sampled from the training set and used as input for the first player, the discriminator, represented by the function D . The goal of the discriminator is to output the probability that its input is real rather than fake, under the assumption that half of the inputs it is ever shown are real and half are fake. In this first scenario, the goal of the discriminator is for $D(\mathbf{x})$ to be near 1. In the second scenario, inputs \mathbf{z} to the generator are randomly sampled from the model's prior over the latent variables. The discriminator then receives input $G(\mathbf{z})$, a fake sample created by the generator. In this scenario, both players participate. The discriminator strives to make $D(G(\mathbf{z}))$ approach 0 while the generative strives to make the same quantity approach 1. If both models have sufficient capacity, then the Nash equilibrium of this game corresponds to the $G(\mathbf{z})$ being drawn from the same distribution as the training data, and $D(\mathbf{x}) = \frac{1}{2}$ for all \mathbf{x} .

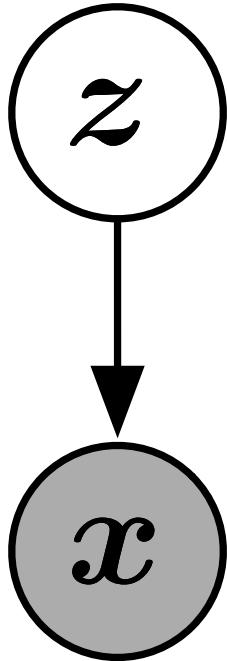


Figure 13: The graphical model structure of GANs, which is also shared with VAEs, sparse coding, etc. It is directed graphical model where every latent variable influences every observed variable. Some GAN variants remove some of these connections.

must be differentiable, but those are the only requirements. In particular, note that any model that can be trained with the nonlinear ICA approach can be a GAN generator network. The relationship with variational autoencoders is more complicated; the GAN framework can train some models that the VAE framework cannot and vice versa, but the two frameworks also have a large intersection. The most salient difference is that, if relying on standard backprop, VAEs cannot have discrete variables at the input to the generator, while GANs cannot have discrete variables at the output of the generator.

The training process The training process consists of simultaneous SGD. On each step, two minibatches are sampled: a minibatch of \mathbf{x} values from the dataset and a minibatch of \mathbf{z} values drawn from the model's prior over latent variables. Then two gradient steps are made simultaneously: one updating $\theta^{(D)}$ to reduce $J^{(D)}$ and one updating $\theta^{(G)}$ to reduce $J^{(G)}$. In both cases, it is possible to use the gradient-based optimization algorithm of your choice. Adam (Kingma and Ba, 2014) is usually a good choice. Many authors recommend running more steps of one player than the other, but as of late 2016, the author's opinion is that the protocol that works the best in practice is simultaneous gradient descent, with one step for each player.

3.2 Cost functions

Several different cost functions may be used within the GANs framework.

3.2.1 The discriminator's cost, $J^{(D)}$

All of the different games designed for GANs so far use the same cost for the discriminator, $J^{(D)}$. They differ only in terms of the cost used for the generator, $J^{(G)}$.

The cost used for the discriminator is:

$$J^{(D)}(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}) = -\frac{1}{2}\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2}\mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z}))). \quad (8)$$

This is just the **standard cross-entropy cost** that is minimized when training a standard binary classifier with a sigmoid output. The only difference is that the classifier is trained on two minibatches of data; one coming from the dataset, where the label is 1 for all examples, and one coming from the generator, where the label is 0 for all examples.

All versions of the GAN game encourage the discriminator to minimize equation 8. In all cases, the discriminator has the same optimal strategy. The reader is now encouraged to complete the exercise in section 7.1 and review its solution given in section 8.1. This exercise shows how to derive the optimal discriminator strategy and discusses the importance of the form of this solution.

We see that by training the discriminator, we are able to obtain an estimate of the ratio

$$\frac{p_{\text{data}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})} \quad (9)$$

at every point \mathbf{x} . Estimating this ratio enables us to compute a wide variety of divergences and their gradients. This is the key approximation technique that sets GANs apart from variational autoencoders and Boltzmann machines. Other deep generative models make approximations based on lower bounds or Markov chains; GANs make approximations based on using supervised learning to estimate a ratio of two densities. The GAN approximation is subject to the failures of supervised learning: overfitting and underfitting. In principle, with perfect optimization and enough training data, these failures can be overcome. Other models make other approximations that have other failures.

Because the GAN framework can naturally be analyzed with the tools of game theory, we call GANs “adversarial.” But we can also think of them as cooperative, in the sense that the discriminator estimates this ratio of densities and then freely shares this information with the generator. From this point of view, the discriminator is more like a teacher instructing the generator in how to improve than an adversary. So far, this cooperative view has not led to any particular change in the development of the mathematics.

3.2.2 Minimax

So far we have specified the cost function for only the discriminator. A complete specification of the game requires that we specify a cost function also for the generator.

The simplest version of the game is a **zero-sum game**, in which the sum of all player's costs is always zero. In this version of the game,

$$J^{(G)} = -J^{(D)}. \quad (10)$$

Because $J^{(G)}$ is tied directly to $J^{(D)}$, we can summarize the entire game with a **value function** specifying the discriminator's payoff:

$$V(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}) = -J^{(D)}(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}). \quad (11)$$

Zero-sum games are also called **minimax** games because their solution involves minimization in an outer loop and maximization in an inner loop:

$$\boldsymbol{\theta}^{(G)*} = \arg \min_{\boldsymbol{\theta}^{(G)}} \max_{\boldsymbol{\theta}^{(D)}} V(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}). \quad (12)$$

The minimax game is mostly of interest because it is easily amenable to theoretical analysis. Goodfellow *et al.* (2014b) used this variant of the GAN game to show that learning in this game resembles minimizing the Jensen-Shannon divergence between the data and the model distribution, and that the game converges to its equilibrium if both players' policies can be updated directly in function space. In practice, the players are represented with deep neural nets and updates are made in parameter space, so these results, which depend on convexity, do not apply.

3.2.3 Heuristic, non-saturating game

The cost used for the generator in the minimax game (equation 10) is useful for theoretical analysis, but does not perform especially well in practice.

Minimizing the cross-entropy between a target class and a classifier's predicted distribution is highly effective because the cost never saturates when the classifier has the wrong output. The cost does eventually saturate, approaching zero, but only when the classifier has already chosen the correct class.

In the minimax game, the discriminator minimizes a cross-entropy, but the generator maximizes the same cross-entropy. This is unfortunate for the generator, because when the discriminator successfully rejects generator samples with high confidence, the generator's gradient vanishes.

To solve this problem, one approach is to continue to use cross-entropy minimization for the generator. Instead of flipping the sign on the discriminator's cost to obtain a cost for the generator, we flip the target used to construct the cross-entropy cost. The cost for the generator then becomes:

$$J^{(G)} = -\frac{1}{2}\mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z})) \quad (13)$$

In the minimax game, the generator minimizes the log-probability of the discriminator being correct. In this game, the generator maximizes the log-probability of the discriminator being mistaken.

This version of the game is heuristically motivated, rather than being motivated by a theoretical concern. The sole motivation for this version of the game is to ensure that each player has a strong gradient when that player is “losing” the game.

In this version of the game, the game is no longer zero-sum, and cannot be described with a single value function.

3.2.4 Maximum likelihood game

We might like to be able to do maximum likelihood learning with GANs, which would mean minimizing the KL divergence between the data and the model, as in equation 4. Indeed, in section 2, we said that GANs could optionally implement maximum likelihood, for the purpose of simplifying their comparison to other models.

There are a variety of methods of approximating equation 4 within the GAN framework. Goodfellow (2014) showed that using

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_z \exp(\sigma^{-1}(D(G(z))), \quad (14)$$

where σ is the logistic sigmoid function, is equivalent to minimizing equation 4, under the assumption that the discriminator is optimal. This equivalence holds in expectation; in practice, both stochastic gradient descent on the KL divergence and the GAN training procedure will have some variance around the true expected gradient due to the use of sampling (of x for maximum likelihood and z for GANs) to construct the estimated gradient. The demonstration of this equivalence is an exercise (section 7.3 with the solution in section 8.3).

Other methods of approximating maximum likelihood within the GANs framework are possible. See for example Nowozin *et al.* (2016).

3.2.5 Is the choice of divergence a distinguishing feature of GANs?

As part of our investigation of how GANs work, we might wonder exactly what it is that makes them work well for generating samples.

Previously, many people (including the author) believed that GANs produced sharp, realistic samples because they minimize the Jensen-Shannon divergence while VAEs produce blurry samples because they minimize the KL divergence between the data and the model.

The KL divergence is not symmetric; minimizing $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$ is different from minimizing $D_{\text{KL}}(p_{\text{model}} \| p_{\text{data}})$. Maximum likelihood estimation performs the former; minimizing the Jensen-Shannon divergence is somewhat more similar to the latter. As shown in figure 14, the latter might be expected to yield better samples because a model trained with this divergence would prefer to generate samples that come only from modes in the training distribution

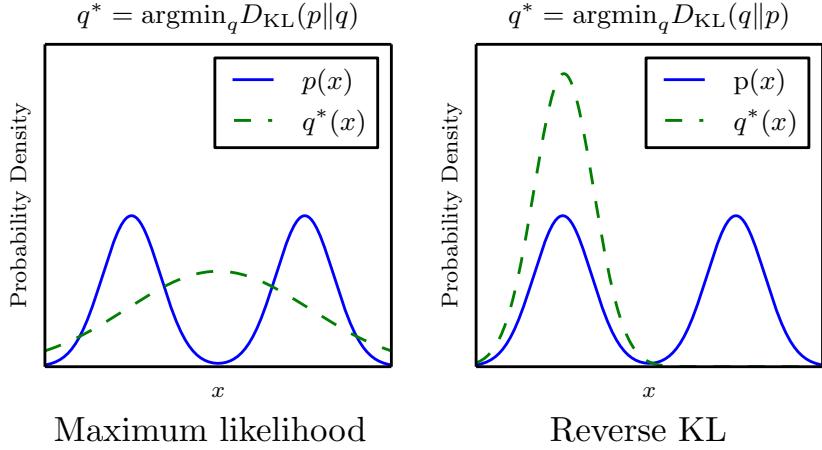


Figure 14: The two directions of the KL divergence are not equivalent. The differences are most obvious when the model has too little capacity to fit the data distribution. Here we show an example of a distribution over one-dimensional data x . In this example, we use a mixture of two Gaussians as the data distribution, and a single Gaussian as the model family. Because a single Gaussian cannot capture the true data distribution, the choice of divergence determines the tradeoff that the model makes. On the left, we use the maximum likelihood criterion. The model chooses to average out the two modes, so that it places high probability on both of them. On the right, we use the reverse order of the arguments to the KL divergence, and the model chooses to capture only one of the two modes. It could also have chosen the other mode; the two are both local minima of the reverse KL divergence. We can think of $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$ as preferring to place high probability everywhere that the data occurs, and $D_{\text{KL}}(p_{\text{model}} \| p_{\text{data}})$ as preferring to place low probability wherever the data does not occur. From this point of view, one might expect $D_{\text{KL}}(p_{\text{model}} \| p_{\text{data}})$ to yield more visually pleasing samples, because the model will not choose to generate unusual samples lying between modes of the data generating distribution.

even if that means ignoring some modes, rather than including all modes but generating some samples that do not come from any training set mode.

Some newer evidence suggests that the use of the Jensen-Shannon divergence does not explain why GANs make sharper samples:

- It is now possible to train GANs using maximum likelihood, as described in section 3.2.4. These models still generate sharp samples, and still select a small number of modes. See figure 15.
- GANs often choose to generate from very few modes; fewer than the limitation imposed by the model capacity. The reverse KL prefers to generate from *as many modes of the data distribution as the model is able to*; it does not prefer fewer modes in general. This suggests that the mode collapse is driven by a factor other than the choice of divergence.

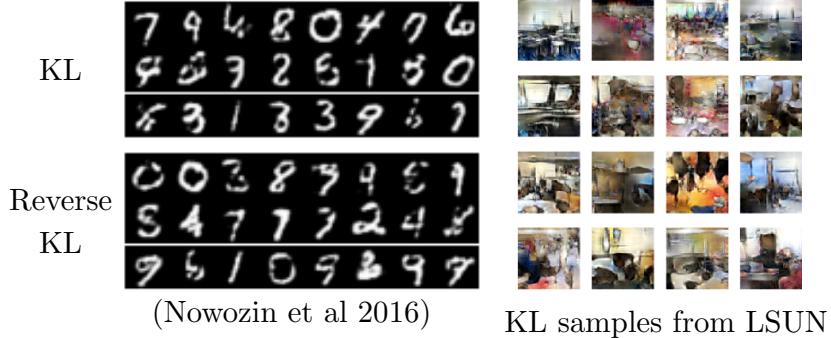


Figure 15: The f-GAN model is able to minimize many different divergences. Because models trained to minimize $D_{\text{KL}}(p_{\text{data}} \parallel p_{\text{model}})$ still generate sharp samples and tend to select a small number of modes, we can conclude that the use of the Jensen-Shannon divergence is not a particular important distinguishing characteristic of GANs and that it does not explain why their samples tend to be sharp.

Altogether, this suggests that GANs choose to generate a small number of modes due to a defect in the training procedure, rather than due to the divergence they aim to minimize. This is discussed further in section 5.1.1. The reason that GANs produce sharp samples is not entirely clear. It may be that the family of models trained using GANs is different from the family of models trained using VAEs (for example, with GANs it is easy to make models where \mathbf{x} has a more complicated distribution than just an isotropic Gaussian conditioned on the input to the generator). It may also be that the kind of approximations that GANs make have different effects than the kind of approximations that other frameworks make.

3.2.6 Comparison of cost functions

We can think of the generator network as learning by a strange kind of reinforcement learning. Rather than being told a specific output \mathbf{x} it should associate with each \mathbf{z} , the generator takes actions and receives rewards for them. In particular, note that $J^{(G)}$ does not make reference to the training data directly at all; all information about the training data comes only through what the discriminator has learned. (Incidentally, this makes GANs resistant to overfitting, because the generator has no opportunity in practice to directly copy training examples) The learning process differs somewhat from traditional reinforcement learning because

- The generator is able to observe not just the output of the reward function but also its gradients.
- The reward function is non-stationary; the reward is based on the discriminator which learns in response to changes in the generator's policy.

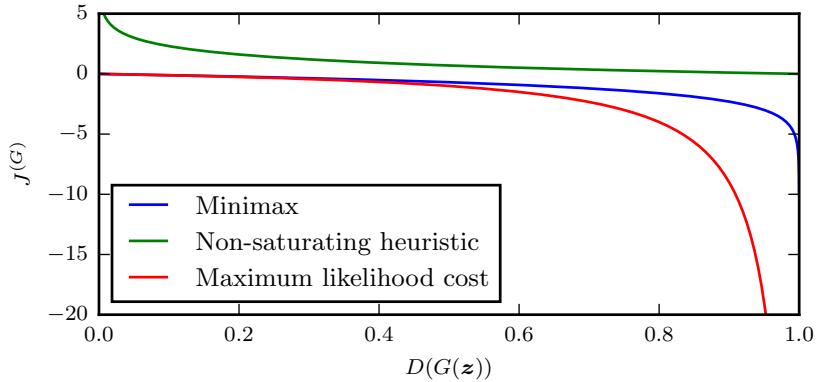


Figure 16: The cost that the generator receives for generating a samples $G(\mathbf{z})$ depends only on how the discriminator responds to that sample. The more probability the discriminator assigns to the sample being real, the less cost the generator receives. We see that when the sample is likely to be fake, both the minimax game and the maximum likelihood game have very little gradient, on the flat left end of the curve. The heuristically motivated non-saturating cost avoids this problem. Maximum likelihood also suffers from the problem that nearly all of the gradient comes from the right end of the curve, meaning that a very small number of samples dominate the gradient computation for each minibatch. This suggests that variance reduction techniques could be an important research area for improving the performance of GANs, especially GANs based on maximum likelihood. Figure reproduced from Goodfellow (2014).

In all cases, we can think of the sampling process that begins with the selection of a specific \mathbf{z} value as an episode that receives a single reward, independent of the actions taken for all other \mathbf{z} values. The reward given to the generator is a function of a single scalar value, $D(G(\mathbf{z}))$. We usually think of this in terms of cost (negative reward). The cost for the generator is always monotonically decreasing in $D(G(\mathbf{z}))$ but different games are designed to make this cost decrease faster along different parts of the curve.

Figure 16 shows the cost response curves as functions of $D(G(\mathbf{z}))$ for three different variants of the GAN game. We see that the maximum likelihood game gives very high variance in the cost, with most of the cost gradient coming from the very few samples of \mathbf{z} that correspond to the samples that are most likely to be real rather than fake. The heuristically designed non-saturating cost has lower sample variance, which may explain why it is more successful in practice. This suggests that [variance reduction techniques](#) could be an important research area for improving the performance of GANs, especially GANs based on maximum likelihood.

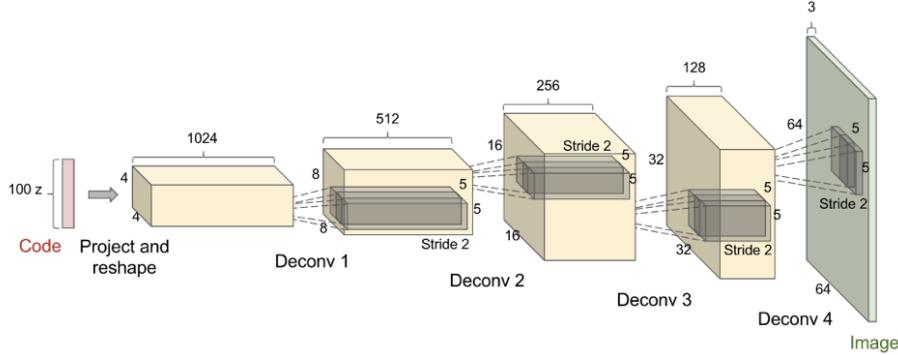


Figure 17: The generator network used by a DCGAN. Figure reproduced from Radford *et al.* (2015).

3.3 The DCGAN architecture

Most GANs today are at least loosely based on the DCGAN architecture (Radford *et al.*, 2015). DCGAN stands for “deep, convolution GAN.” Though GANs were both deep and convolutional prior to DCGANs, the name DCGAN is useful to refer to this specific style of architecture. Some of the key insights of the DCGAN architecture were to:

- Use batch normalization (Ioffe and Szegedy, 2015) layers in most layers of both the discriminator and the generator, with the two minibatches for the discriminator normalized separately. The last layer of the generator and first layer of the discriminator are not batch normalized, so that the model can learn the correct mean and scale of the data distribution. See figure 17.
- The overall network structure is mostly borrowed from the all-convolutional net (Springenberg *et al.*, 2015). This architecture contains neither pooling nor “unpooling” layers. When the generator needs to increase the spatial dimension of the representation it uses transposed convolution with a stride greater than 1.
- The use of the Adam optimizer rather than SGD with momentum.

Prior to DCGANs, LAPGANs (Denton *et al.*, 2015) were the only version of GAN that had been able to scale to high resolution images. LAPGANs require a multi-stage generation process in which multiple GANs generate different levels of detail in a Laplacian pyramid representation of an image. DCGANs were the first GAN model to learn to generate high resolution images in a single shot. As shown in figure 18, DCGANs are able to generate high quality images when trained on restricted domains of images, such as images of bedrooms. DCGANs also clearly demonstrated that GANs learn to use their latent code in meaningful ways, with simple arithmetic operations in latent space having clear

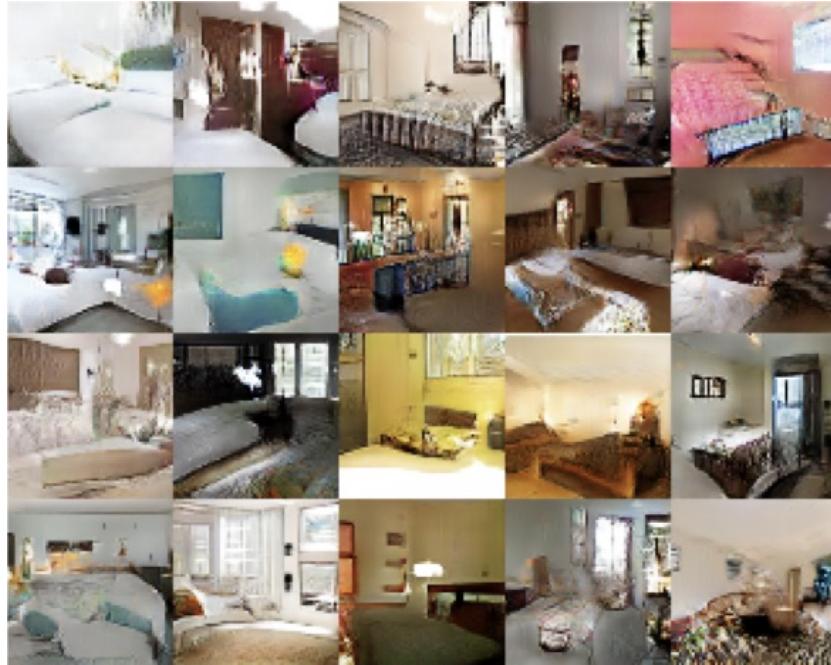


Figure 18: Samples of images of bedrooms generated by a DCGAN trained on the LSUN dataset.

interpretation as arithmetic operations on semantic attributes of the input, as demonstrated in figure 19.

3.4 How do GANs relate to noise-contrastive estimation and maximum likelihood?

While trying to understand how GANs work, one might naturally wonder about how they are connected to **noise-contrastive estimation** (NCE) (Gutmann and Hyvarinen, 2010). Minimax GANs use the cost function from NCE as their value function, so the methods seem closely related at face value. It turns out that they actually learn very different things, because the two methods focus on different players within this game. Roughly speaking, the goal of NCE is to learn the density model within the discriminator, while the goal of GANs is to learn the sampler defining the generator. While these two tasks seem closely related at a qualitative level, the gradients for the tasks are actually quite different. Surprisingly, maximum likelihood turns out to be closely related to NCE, and corresponds to playing a minimax game with the same value function, but using a sort of heuristic update strategy rather than gradient descent for one of the two players. The connections are summarized in figure 20.

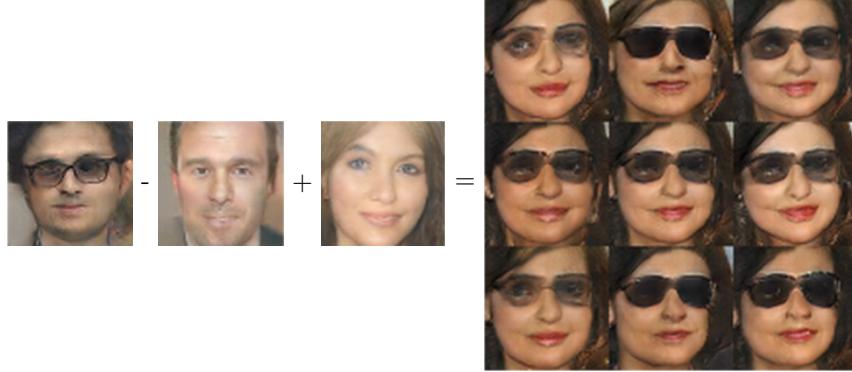


Figure 19: DCGANs demonstrated that GANs can learn a distributed representation that disentangles the concept of gender from the concept of wearing glasses. If we begin with the representation of the concept of a man with glasses, then subtract the vector representing the concept of a man without glasses, and finally add the vector representing the concept of a woman without glasses, we obtain the vector representing the concept of a woman with glasses. The generative model correctly decodes all of these representation vectors to images that may be recognized as belonging to the correct class. Images reproduced from Radford *et al.* (2015).

$$V(G, D) = \mathbb{E}_{p_{\text{data}}} \log D(\mathbf{x}) + \mathbb{E}_{p_{\text{generator}}} (\log (1 - D(\mathbf{x})))$$

	NCE (Gutmann and Hyvärinen 2010)	MLE	GAN
D	$D(x) = \frac{p_{\text{model}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x}) + p_{\text{generator}}(\mathbf{x})}$		Neural network
Goal	Learn p_{model}	Learn $p_{\text{generator}}$	
G update rule	None (G is fixed)	Copy p_{model} parameters	Gradient descent on V
D update rule	Gradient ascent on V		

Figure 20: Goodfellow (2014) demonstrated the following connections between minimax GANs, noise-contrastive estimation, and maximum likelihood: all three can be interpreted as strategies for playing a minimax game with the same value function. The biggest difference is in where p_{model} lies. For GANs, the generator is p_{model} , while for NCE and MLE, p_{model} is part of the discriminator. Beyond this, the differences between the methods lie in the update strategy. GANs learn both players with gradient descent. MLE learns the discriminator using gradient descent, but has a heuristic update rule for the generator. Specifically, after each discriminator update step, MLE copies the density model learned inside the discriminator and converts it into a sampler to be used as the generator. NCE never updates the generator; it is just a fixed source of noise.

4 Tips and Tricks

Practitioners use several tricks to improve the performance of GANs. It can be difficult to tell how effective some of these tricks are; many of them seem to help in some contexts and hurt in others. These should be regarded as techniques that are worth trying out, not as ironclad best practices.

NIPS 2016 also featured a workshop on adversarial training, with an invited talk by Soumith Chintala called "How to train a GAN." This talk has more or less the same goal as this portion of the tutorial, with a different collection of advice. To learn about tips and tricks not included in this tutorial, check out the GitHub repository associated with Soumith's talk:

<https://github.com/soumith/ganhacks>

4.1 Train with labels

Using labels in any way, shape or form almost always results in a dramatic improvement in the subjective quality of the samples generated by the model. This was first observed by Denton *et al.* (2015), who built class-conditional GANs that generated much better samples than GANs that were free to generate from any class. Later, Salimans *et al.* (2016) found that sample quality improved even if the generator did not explicitly incorporate class information; training the discriminator to recognize specific classes of real objects is sufficient.

It is not entirely clear why this trick works. It may be that the incorporation of class information gives the training process useful clues that help with optimization. It may also be that this trick gives no objective improvement in sample quality, but instead biases the samples toward taking on properties that the human visual system focuses on. If the latter is the case, then this trick may not result in a better model of the true data-generating distribution, but it still helps to create media for a human audience to enjoy and may help an RL agent to carry out tasks that rely on knowledge of the same aspects of the environment that are relevant to human beings.

It is important to compare results obtained using this trick only to other results using the same trick; models trained with labels should be compared only to other models trained with labels, class-conditional models should be compared only to other class-conditional models. Comparing a model that uses labels to one that does not is unfair and an uninteresting benchmark, much as a convolutional model can usually be expected to outperform a non-convolutional model on image tasks.

4.2 One-sided label smoothing

GANs are intended to work when the discriminator estimates a ratio of two densities, but deep neural nets are prone to producing highly confident outputs that identify the correct class but with too extreme of a probability. This is especially the case when the input to the deep network is adversarially constructed;

the classifier tends to linearly extrapolate and produce extremely confident predictions (Goodfellow *et al.*, 2014a).

To encourage the discriminator to estimate soft probabilities rather than to extrapolate to extremely confident classification, we can use a technique called **one-sided label smoothing** (Salimans *et al.*, 2016).

Usually we train the discriminator using equation 8. We can write this in TensorFlow (Abadi *et al.*, 2015) code as:

```
d_on_data = discriminator_logits(data_minibatch)
d_on_samples = discriminator_logits(samples_minibatch)
loss = tf.nn.sigmoid_cross_entropy_with_logits(d_on_data, 1.) + \
       tf.nn.sigmoid_cross_entropy_with_logits(d_on_samples, 0.)
```

The idea of one-sided label smoothing is to replace the target for the real examples with a value slightly less than one, such as .9:

```
loss = tf.nn.sigmoid_cross_entropy_with_logits(d_on_data, .9) + \
       tf.nn.sigmoid_cross_entropy_with_logits(d_on_samples, 0.)
```

This prevents extreme extrapolation behavior in the discriminator; if it learns to predict extremely large logits corresponding to a probability approaching 1 for some input, it will be penalized and encouraged to bring the logits back down to a smaller value.

It is important to not smooth the labels for the fake samples. Suppose we use a target of $1 - \alpha$ for the real data and a target of $0 + \beta$ for the fake samples. Then the optimal discriminator function is

$$D^*(\mathbf{x}) = \frac{(1 - \alpha)p_{\text{data}}(\mathbf{x}) + \beta p_{\text{model}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}. \quad (15)$$

When β is zero, then smoothing by α does nothing but scale down the optimal value of the discriminator. When β is nonzero, the shape of the optimal discriminator function changes. In particular, in a region where $p_{\text{data}}(\mathbf{x})$ is very small and $p_{\text{model}}(\mathbf{x})$ is larger, $D^*(\mathbf{x})$ will have a peak near the spurious mode of $p_{\text{model}}(\mathbf{x})$. The discriminator will thus reinforce incorrect behavior in the generator; the generator will be trained either to produce samples that resemble the data or to produce samples that resemble the samples it already makes.

One-sided label smoothing is a simple modification of the much older label smoothing technique, which dates back to at least the 1980s. Szegedy *et al.* (2015) demonstrated that label smoothing is an excellent regularizer in the context of convolutional networks for object recognition. One reason that label smoothing works so well as a regularizer is that it does not ever encourage the model to choose an incorrect class on the training set, but only to reduce the confidence in the correct class. Other regularizers such as weight decay often encourage some misclassification if the coefficient on the regularizer is set high enough. Warde-Farley and Goodfellow (2016) showed that label smoothing can help to reduce vulnerability to adversarial examples, which suggests that label smoothing should help the discriminator more efficiently learn to resist attack by the generator.

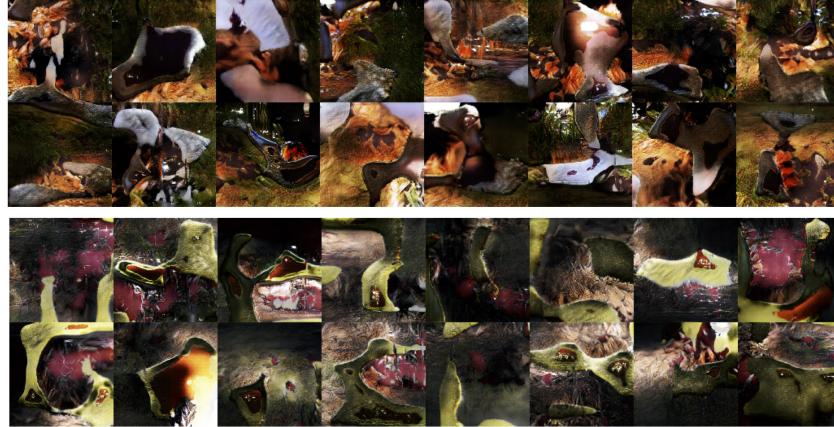


Figure 21: Two minibatches of sixteen samples each, generated by a generator network using batch normalization. These minibatches illustrate a problem that occurs occasionally when using batch normalization: **fluctuations in the mean and standard deviation of feature values in a minibatch can have a greater effect than the individual z codes for individual images within the minibatch.** This manifests here as one minibatch containing all orange-tinted samples and the other containing all green-tinted samples. The examples within a minibatch should be independent from each other, but in this case, batch normalization has caused them to become correlated with each other.

4.3 Virtual batch normalization

Since the introduction of DCGANs, most GAN architectures have involved some form of batch normalization. The main purpose of batch normalization is to improve the optimization of the model, by reparameterizing the model so that the mean and variance of each feature are determined by a single mean parameter and a single variance parameter associated with that feature, rather than by a complicated interaction between all of the weights of all of the layers used to extract the feature. This reparameterization is accomplished by subtracting the mean and dividing by the standard deviation of that feature on a minibatch of data. It is important that the normalization operation is *part of the model*, so that back-propagation computes the gradient of features that are defined to always be normalized. The method is much less effective if features are frequently renormalized after learning without the normalization defined as part of the model.

Batch normalization is very helpful, but for GANs has a few unfortunate side effects. The use of a different minibatch of data to compute the normalization statistics on each step of training results in fluctuation of these normalizing constants. When minibatch sizes are small (as is often the case when trying to fit a large generative model into limited GPU memory) these fluctuations can become large enough that they have more effect on the image generated by the GAN than the input z has. See figure 21 for an example.

Salimans *et al.* (2016) introduced techniques to mitigate this problem. **Reference batch normalization** consists of running the network twice: once on a minibatch of **reference examples** that are sampled once at the start of training and never replaced, and once on the current minibatch of examples to train on. The mean and standard deviation of each feature are computed using the reference batch. The features for both batches are then normalized using these computed statistics. A drawback to reference batch normalization is that the model can overfit to the reference batch. To mitigate this problem slightly, one can instead use **virtual batch normalization**, in which the normalization statistics for each example are computed using the union of that example and the reference batch. Both reference batch normalization and virtual batch normalization have the property that all examples in the training minibatch are processed independently from each other, and all samples produced by the generator (except those defining the reference batch) are i.i.d.

4.4 Can one balance G and D ?

Many people have an intuition that it is necessary to somehow balance the two players to prevent one from overpowering the other. If such balance is desirable and feasible, it has not yet been demonstrated in any compelling fashion.

The author’s present belief is that GANs work by estimating the ratio of the data density and model density. This ratio is estimated correctly only when the discriminator is optimal, so it is fine for the discriminator to overpower the generator.

Sometimes the gradient for the generator can vanish when the discriminator becomes too accurate. The right way to solve this problem is not to limit the power of the discriminator, but to use a parameterization of the game where the gradient does not vanish (section 3.2.3).

Sometimes the gradient for the generator can become very large if the discriminator becomes too confident. Rather than making the discriminator less accurate, a better way to resolve this problem is to use one-sided label smoothing (section 4.2).

The idea that the discriminator should always be optimal in order to best estimate the ratio would suggest training the discriminator for $k > 1$ steps every time the generator is trained for one step. In practice, this does not usually result in a clear improvement.

One can also try to balance the generator and discriminator by choosing the model size. In practice, the discriminator is usually deeper and sometimes has more filters per layer than the generator. This may be because it is important for the discriminator to be able to correctly estimate the ratio between the data density and generator density, but it may also be an artifact of the mode collapse problem—since the generator tends not to use its full capacity with current training methods, practitioners presumably do not see much of a benefit from increasing the generator capacity. If the mode collapse problem can be overcome, generator sizes will presumably increase. It is not clear whether discriminator sizes will increase proportionally.

5 Research Frontiers

GANs are a relatively new method, with many research directions still remaining open.

5.1 Non-convergence

The largest problem facing GANs that researchers should try to resolve is the issue of non-convergence.

Most deep models are trained using an optimization algorithm that seeks out a low value of a cost function. While many problems can interfere with optimization, optimization algorithms usually make reliable downhill progress. GANs require finding the equilibrium to a game with two players. Even if each player successfully moves downhill on that player’s update, the same update might move the other player uphill. Sometimes the two players eventually reach an equilibrium, but in other scenarios they repeatedly undo each others’ progress without arriving anywhere useful. This is a general problem with games not unique to GANs, so a general solution to this problem would have wide-reaching applications.

To gain some intuition for how gradient descent performs when applied to games rather than optimization, the reader is encouraged to solve the exercise in section 7.2 and review its solution in section 8.2 now.

Simultaneous gradient descent converges for some games but not all of them.

In the case of the minimax GAN game (section 3.2.2), Goodfellow *et al.* (2014b) showed that simultaneous gradient descent converges *if the updates are made in function space*. In practice, the updates are made in parameter space, so the convexity properties that the proof relies on do not apply. Currently, there is neither a theoretical argument that GAN games should converge when the updates are made to parameters of deep neural networks, nor a theoretical argument that the games should not converge.

In practice, GANs often seem to oscillate, somewhat like what happens in the toy example in section 8.2, meaning that they progress from generating one kind of sample to generating another kind of sample without eventually reaching an equilibrium.

Probably the most common form of harmful non-convergence encountered in the GAN game is mode collapse.

5.1.1 Mode collapse

Mode collapse, also known as **the Helvetica scenario**, is a problem that occurs when the generator learns to map several different input z values to the same output point. In practice, complete mode collapse is rare, but partial mode collapse is common. **Partial mode collapse** refers to scenarios in which the generator makes multiple images that contain the same color or texture themes, or multiple images containing different views of the same dog. The mode collapse problem is illustrated in figure 22.

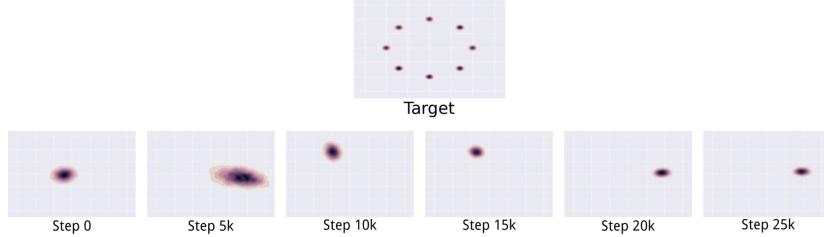


Figure 22: An illustration of the mode collapse problem on a two-dimensional toy dataset. In the top row, we see the target distribution p_{data} that the model should learn. It is a mixture of Gaussians in a two-dimensional space. In the lower row, we see a series of different distributions learned over time as the GAN is trained. Rather than converging to a distribution containing all of the modes in the training set, the generator only ever produces a single mode at a time, cycling between different modes as the discriminator learns to reject each one. Images from Metz *et al.* (2016).

Mode collapse may arise because the maximin solution to the GAN game is different from the minimax solution. When we find the model

$$G^* = \min_G \max_D V(G, D), \quad (16)$$

G^* draws samples from the data distribution. When we exchange the order of the min and max and find

$$G^* = \max_D \min_G V(G, D), \quad (17)$$

the minimization with respect to the generator now lies in the inner loop of the optimization procedure. The generator is thus asked to map every \mathbf{z} value to the single \mathbf{x} coordinate that the discriminator believes is most likely to be real rather than fake. Simultaneous gradient descent does not clearly privilege min max over max min or vice versa. We use it in the hope that it will behave like min max but it often behaves like max min.

As discussed in section 3.2.5, mode collapse does not seem to be caused by any particular cost function. It is commonly asserted that mode collapse is caused by the use of Jensen-Shannon divergence, but this does not seem to be the case, because GANs that minimize approximations of $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$ face the same issues, and because the generator often collapses to even fewer modes than would be preferred by the Jensen-Shannon divergence.

Because of the mode collapse problem, applications of GANs are often limited to problems where it is acceptable for the model to produce a small number of distinct outputs, usually tasks where the goal is to map some input to one of many acceptable outputs. As long as the GAN is able to find a small number of these acceptable outputs, it is useful. One example is text-to-image synthesis, in which the input is a caption for an image, and the output is an image matching that description. See figure 23 for a demonstration of this task. In very recent work, Reed *et al.* (2016a) have shown that other models have higher

this small bird has a pink breast and crown, and black primaries and secondaries. this magnificent fellow is almost all black with a red crest, and white cheek patch.



the flower has petals that are bright pinkish purple with white stigma



this white and yellow flower have thin white petals and a round yellow stamen



Figure 23: Text-to-image synthesis with GANs. Image reproduced from Reed *et al.* (2016b).

output diversity than GANs for such tasks (figure 24), but StackGANs (Zhang *et al.*, 2016) seem to have higher output diversity than previous GAN-based approaches (figure 25).

The mode collapse problem is probably the most important issue with GANs that researchers should attempt to address.

One attempt is **minibatch features** (Salimans *et al.*, 2016). The basic idea of minibatch features is to allow the discriminator to compare an example to a minibatch of generated samples and a minibatch of real samples. By measuring distances to these other samples in latent spaces, the discriminator can detect if a sample is unusually similar to other generated samples. Minibatch features work well. It is strongly recommended to directly copy the Theano/TensorFlow code released with the paper that introduced them, since small changes in the definition of the features result in large reductions in performance.

Minibatch GANs trained on CIFAR-10 obtain excellent results, with most samples being recognizable as specific CIFAR-10 classes (figure 26). When



Figure 24: GANs have low output diversity for text-to-image tasks because of the mode collapse problem. Image reproduced from Reed *et al.* (2016a).

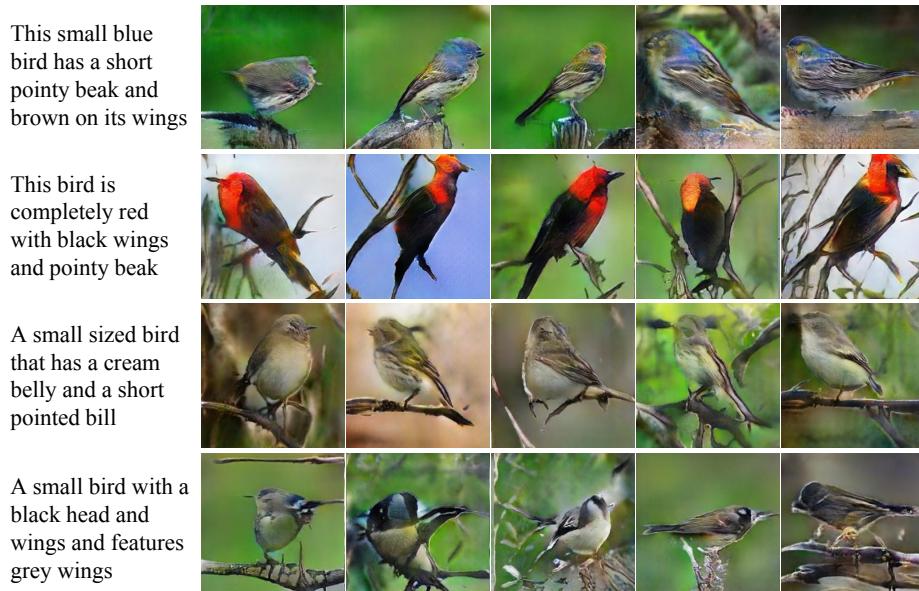


Figure 25: StackGANs are able to achieve higher output diversity than other GAN-based text-to-image models. Image reproduced from Zhang *et al.* (2016).

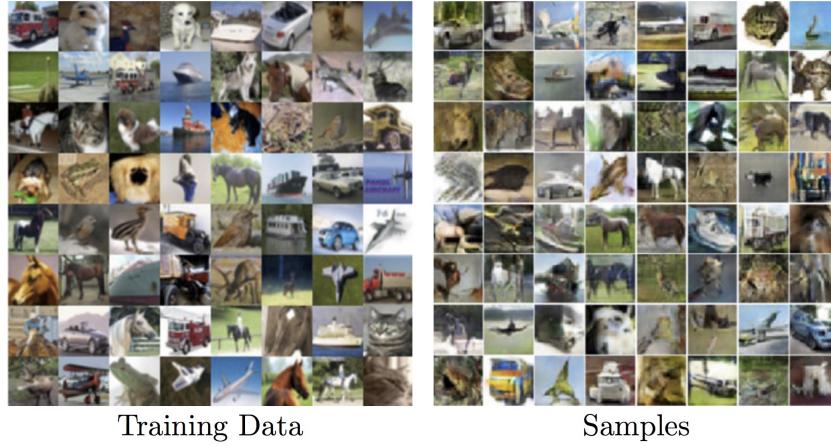


Figure 26: Minibatch GANs trained on CIFAR-10 obtain excellent results, with most samples being recognizable as specific CIFAR-10 classes. (Note: this model was trained with labels)

trained on 128×128 ImageNet, few images are recognizable as belonging to a specific ImageNet class (figure 27). Some of the better images are cherry-picked into figure 28.

Minibatch GANs have reduced the mode collapse problem enough that other problems, such as difficulties with counting, perspective, and global structure become the most obvious defects (figure 29, figure 30, and figure 31, respectively). Many of these problems could presumably be resolved by designing better model architectures.

Another approach to solving the mode collapse problem is **unrolled GANs** (Metz *et al.*, 2016). Ideally, we would like to find $G^* = \arg \min_G \max_D V(G, D)$. In practice, when we simultaneously follow the gradient of $V(G, D)$ for both players, we essentially ignore the max operation when computing the gradient for G . Really, we should regard $\max_D V(G, D)$ as the cost function for G , and we should back-propagate through the maximization operation. Various strategies exist for back-propagating through a maximization operation, but many, such as those based on implicit differentiation, are unstable. The idea of unrolled GANs is to build a computational graph describing k steps of learning in the discriminator, then backpropagate through all k of these steps of learning when computing the gradient on the generator. Fully maximizing the value function for the discriminator takes tens of thousands of steps, but Metz *et al.* (2016) found that unrolling for even small numbers of steps, like 10 or fewer, can noticeably reduce the mode dropping problem. This approach has not yet been scaled up to ImageNet. See figure 32 for a demonstration of unrolled GANs on a toy problem.

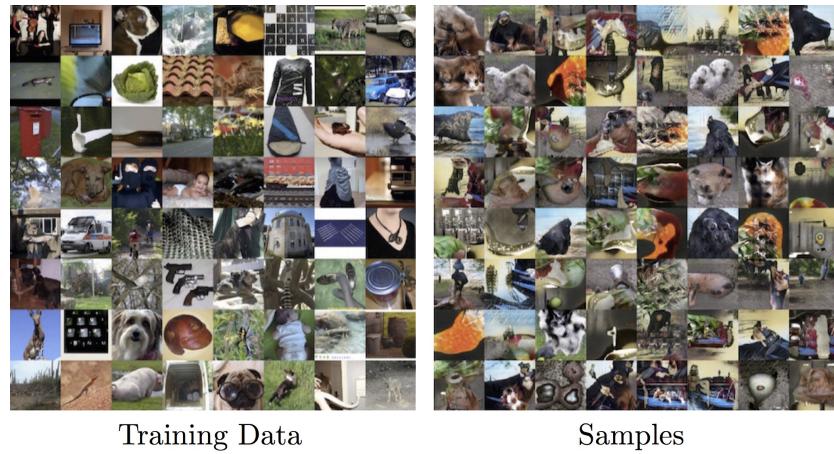


Figure 27: Minibatch GANS trained with labels on 128×128 ImageNet produce images that are occasionally recognizable as belonging to specific classes.

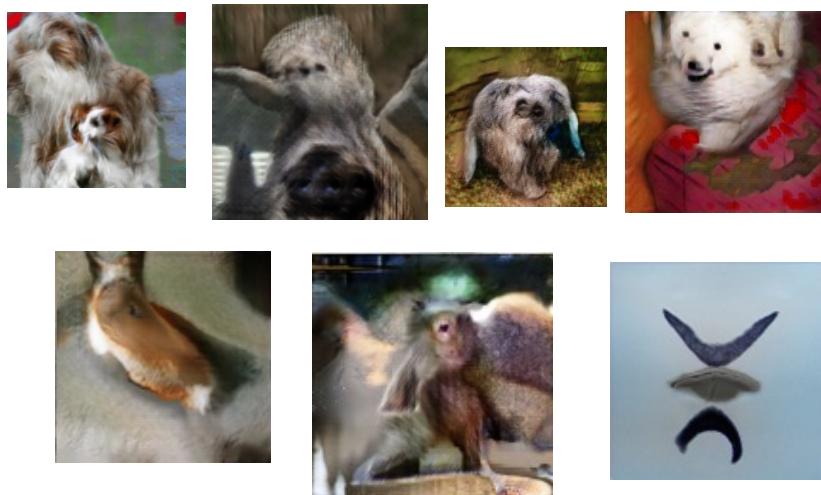


Figure 28: Minibatch GANs sometimes produce very good images when trained on 128×128 ImageNet, as demonstrated by these cherry-picked examples.

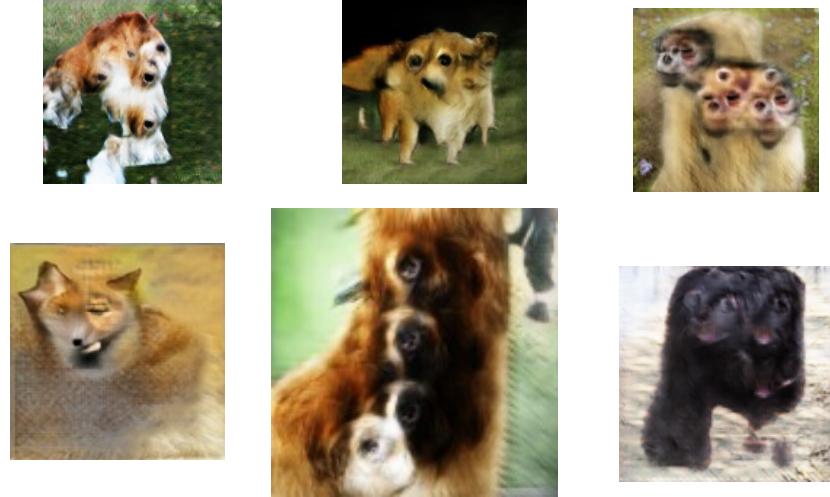


Figure 29: GANs on 128×128 ImageNet seem to have trouble with counting, often generating animals with the wrong number of body parts.



Figure 30: GANs on 128×128 ImageNet seem to have trouble with the idea of three-dimensional perspective, often generating images of objects that are too flat or highly axis-aligned. As a test of the reader's discriminator network, one of these images is actually real.



Figure 31: GANs on 128×128 ImageNet seem to have trouble coordinating global structure, for example, drawing “Fallout Cow,” an animal that has both quadrupedal and bipedal structure.

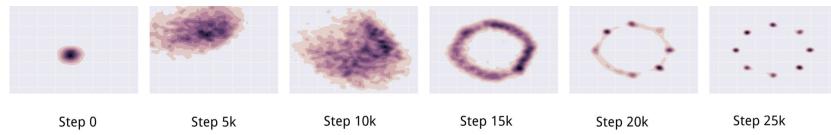


Figure 32: Unrolled GANs are able to fit all of the modes of a mixture of Gaussians in a two-dimensional space. Image reproduced from Metz *et al.* (2016).

5.1.2 Other games

If our theory of how to understand whether a continuous, high-dimensional non-convex game will converge could be improved, or if we could develop algorithms that converge more reliably than simultaneous gradient descent, several application areas besides GANs would benefit. Even restricted to just AI research, we find games in many scenarios:

- Agents that literally play games, such as AlphaGo (Silver *et al.*, 2016).
- Machine learning security, where models must resist adversarial examples (Szegedy *et al.*, 2014; Goodfellow *et al.*, 2014a).
- Domain adaptation via domain-adversarial learning (Ganin *et al.*, 2015).
- Adversarial mechanisms for preserving privacy (Edwards and Storkey, 2015).
- Adversarial mechanisms for cryptography (Abadi and Andersen, 2016).

This is by no means an exhaustive list.

5.2 Evaluation of generative models

Another highly important research area related to GANs is that it is not clear how to quantitatively evaluate generative models. Models that obtain good likelihood can generate bad samples, and models that generate good samples can have poor likelihood. There is no clearly justified way to quantitatively score samples. GANs are somewhat harder to evaluate than other generative models because it can be difficult to estimate the likelihood for GANs (but it is possible—see Wu *et al.* (2016)). Theis *et al.* (2015) describe many of the difficulties with evaluating generative models.

5.3 Discrete outputs

The only real requirement imposed on the design of the generator by the GAN framework is that the generator must be differentiable. Unfortunately, this means that the generator cannot produce discrete data, such as one-hot word or character representations. Removing this limitation is an important research direction that could unlock the potential of GANs for NLP. There are at least three obvious ways one could attack this problem:

1. Using the REINFORCE algorithm (Williams, 1992).
2. Using the concrete distribution (Maddison *et al.*, 2016) or Gumbel-softmax (Jang *et al.*, 2016).
3. Training the generator to sample continuous values that can be decoded to discrete ones (e.g., sampling word embeddings directly).

5.4 Semi-supervised learning

A research area where GANs are already highly successful is the use of generative models for semi-supervised learning, as proposed but not demonstrated in the original GAN paper (Goodfellow *et al.*, 2014b).

GANs have been successfully applied to semi-supervised learning at least since the introduction of CatGANs (Springenberg, 2015). Currently, the state of the art in semi-supervised learning on MNIST, SVHN, and CIFAR-10 is obtained by **feature matching GANs** (Salimans *et al.*, 2016). Typically, models are trained on these datasets using 50,000 or more labels, but feature matching GANs are able to obtain good performance with very few labels. They obtain state of the art performance within several categories for different amounts of labels, ranging from 20 to 8,000.

The basic idea of how to do semi-supervised learning with feature matching GANs is to turn a classification problem with n classes into a classification problem with $n + 1$ classes, with the additional class corresponding to fake images. All of the real classes can be summed together to obtain the probability of the image being real, enabling the use of the classifier as a discriminator within the GAN game. The real-vs-fake discriminator can be trained even with unlabeled data, which is known to be real, and with samples from the generator, which are known to be fake. The classifier can also be trained to recognize individual real classes on the limited amount of real, labeled examples. This approach was simultaneously developed by Salimans *et al.* (2016) and Odena (2016). The earlier CatGAN used an n class discriminator rather than an $n + 1$ class discriminator.

Future improvements to GANs can presumably be expected to yield further improvements to semi-supervised learning.

5.5 Using the code

GANs learn a representation \mathbf{z} of the image \mathbf{x} . It is already known that this representation can capture useful high-level abstract semantic properties of \mathbf{x} , but it can be somewhat difficult to make use of this information.

One obstacle to using \mathbf{z} is that it can be difficult to obtain \mathbf{z} given an input \mathbf{x} . Goodfellow *et al.* (2014b) proposed but did not demonstrate using a second network analogous to the generator to sample from $p(\mathbf{z} | \mathbf{x})$, much as the generator samples from $p(\mathbf{x})$. So far the full version of this idea, using a fully general neural network as the encoder and sampling from an arbitrarily powerful approximation of $p(\mathbf{z} | \mathbf{x})$, has not been successfully demonstrated, but Donahue *et al.* (2016) demonstrated how to train a deterministic encoder, and Dumoulin *et al.* (2016) demonstrated how to train an encoder network that samples from a Gaussian approximation of the posterior. Further research will presumably develop more powerful stochastic encoders.

Another way to make better use of the code is to train the code to be more useful. InfoGANs (Chen *et al.*, 2016a) regularize some entries in the code vector with an extra objective function that encourages them to have high mutual

information with \mathbf{x} . Individual entries in the resulting code then correspond to specific semantic attributes of \mathbf{x} , such as the direction of lighting on an image of a face.

5.6 Developing connections to reinforcement learning

Researchers have already identified connections between GANs and actor-critic methods (Pfau and Vinyals, 2016), inverse reinforcement learning (Finn *et al.*, 2016a), and have applied GANs to imitation learning (Ho and Ermon, 2016). These connections to RL will presumably continue to bear fruit, both for GANs and for RL.

6 Plug and Play Generative Networks

Shortly before this tutorial was presented at NIPS, a new generative model was released. This model, plug and play generative networks (Nguyen *et al.*, 2016), has dramatically improved the diversity of samples of images of ImageNet classes that can be produced at high resolution.

PPGNs are new and not yet well understood. The model is complicated, and most of the recommendations about how to design the model are based on empirical observation rather than theoretical understanding. This tutorial will thus not say too much about exactly how PPGNs work, since this will presumably become more clear in the future.

As a brief summary, PPGNs are basically an approximate Langevin sampling approach to generating images with a Markov chain. The gradients for the Langevin sampler are estimated using a denoising autoencoder. The denoising autoencoder is trained with several losses, including a GAN loss.

Some of the results are shown in figure 33. As demonstrated in figure 34, the GAN loss is crucial for obtaining high quality images.

7 Exercises

This tutorial includes three exercises to check your understanding. The solutions are given in section 8.

7.1 The optimal discriminator strategy

As described in equation 8, the goal of the discriminator is to minimize

$$J^{(D)}(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}) = -\frac{1}{2}\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2}\mathbb{E}_{\mathbf{z}} \log (1 - D(G(z))) \quad (18)$$

with respect to $\boldsymbol{\theta}^{(D)}$. Imagine that the discriminator can be optimized in function space, so the value of $D(\mathbf{x})$ is specified independently for every value of \mathbf{x} . What is the optimal strategy for D ? What assumptions need to be made to obtain this result?

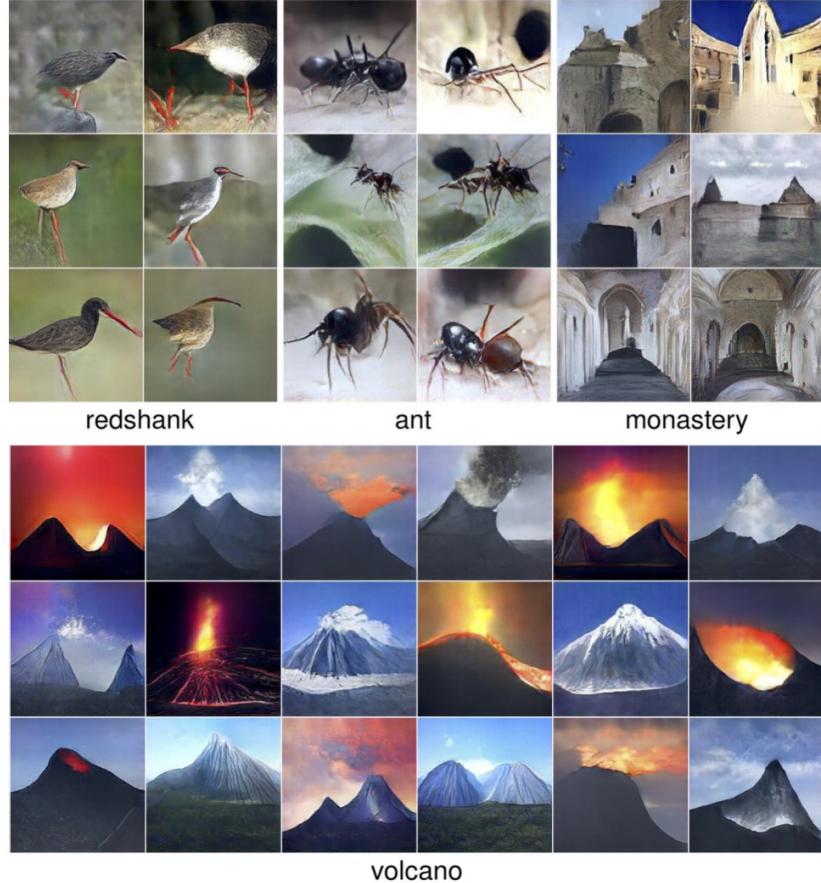


Figure 33: PPGNs are able to generate diverse, high resolution images from ImageNet classes. Image reproduced from Nguyen *et al.* (2016).

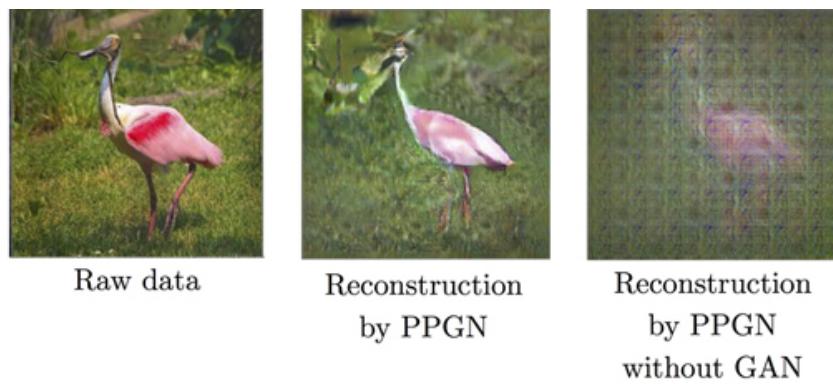


Figure 34: The GAN loss is a crucial ingredient of PPGNs. Without it, the denoising autoencoder used to drive PPGNs does not create compelling images.

7.2 Gradient descent for games

Consider a minimax game with two players that each control a single scalar value. The minimizing player controls scalar x and the maximizing player controls scalar y . The value function for this game is

$$V(x, y) = xy. \quad (19)$$

- Does this game have an equilibrium? If so, where is it?
- Consider the learning dynamics of simultaneous gradient descent. To simplify the problem, treat gradient descent as a continuous time process. With an infinitesimal learning rate, gradient descent is described by a system of partial differential equations:

$$\frac{\partial x}{\partial t} = -\frac{\partial}{\partial x} V(x(t), y(t)) \quad (20)$$

$$\frac{\partial y}{\partial t} = \frac{\partial}{\partial y} V(x(t), y(t)). \quad (21)$$

Solve for the trajectory followed by these dynamics.

7.3 Maximum likelihood in the GAN framework

In this exercise, we will derive a cost that yields (approximate) maximum likelihood learning within the GAN framework. Our goal is to design $J^{(G)}$ so that, if we assume the discriminator is optimal, the expected gradient of $J^{(G)}$ will match the expected gradient of $D_{\text{KL}}(p_{\text{data}} \| p_{\text{model}})$.

The solution will take the form of:

$$J^{(G)} = \mathbb{E}_{\mathbf{x} \sim p_g} f(\mathbf{x}). \quad (22)$$

The exercise consists of determining the form of f .

8 Solutions to exercises

8.1 The optimal discriminator strategy

Our goal is to minimize

$$J^{(D)}(\boldsymbol{\theta}^{(D)}, \boldsymbol{\theta}^{(G)}) = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(z))) \quad (23)$$

in function space, specifying $D(\mathbf{x})$ directly.

We begin by assuming that both p_{data} and p_{model} are nonzero everywhere. If we do not make this assumption, then some points are never visited during training, and have undefined behavior.

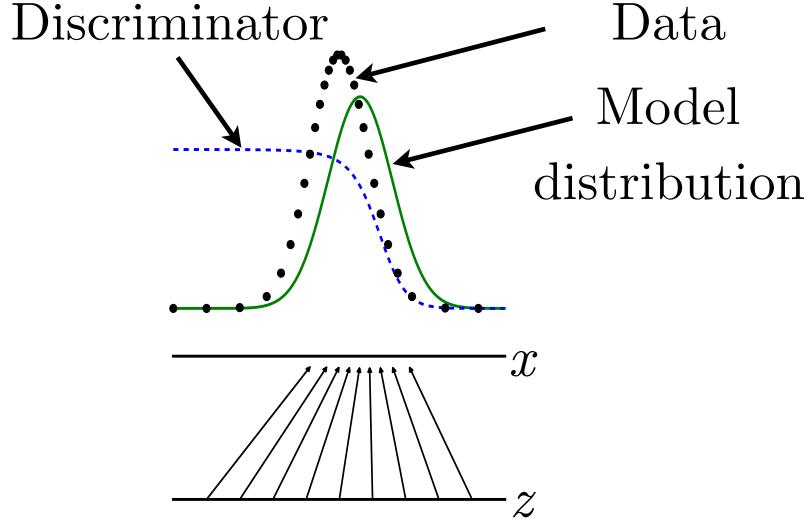


Figure 35: An illustration of how the discriminator estimates a ratio of densities. In this example, we assume that both z and x are one dimensional for simplicity. The mapping from z to x (shown by the black arrows) is non-uniform so that $p_{\text{model}}(x)$ (shown by the green curve) is greater in places where z values are brought together more densely. The discriminator (dashed blue line) estimates the ratio between the data density (black dots) and the sum of the data and model densities. Wherever the output of the discriminator is large, the model density is too low, and wherever the output of the discriminator is small, the model density is too high. The generator can learn to produce a better model density by following the discriminator uphill; each $G(z)$ value should move slightly in the direction that increases $D(G(z))$. Figure reproduced from Goodfellow *et al.* (2014b).

To minimize $J^{(D)}$ with respect to D , we can write down the functional derivatives with respect to a single entry $D(\mathbf{x})$, and set them equal to zero:

$$\frac{\delta}{\delta D(\mathbf{x})} J^{(D)} = 0. \quad (24)$$

By solving this equation, we obtain

$$D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}. \quad (25)$$

Estimating this ratio is the key approximation mechanism used by GANs. The process is illustrated in figure 35.

8.2 Gradient descent for games

The value function

$$V(x, y) = xy \quad (26)$$

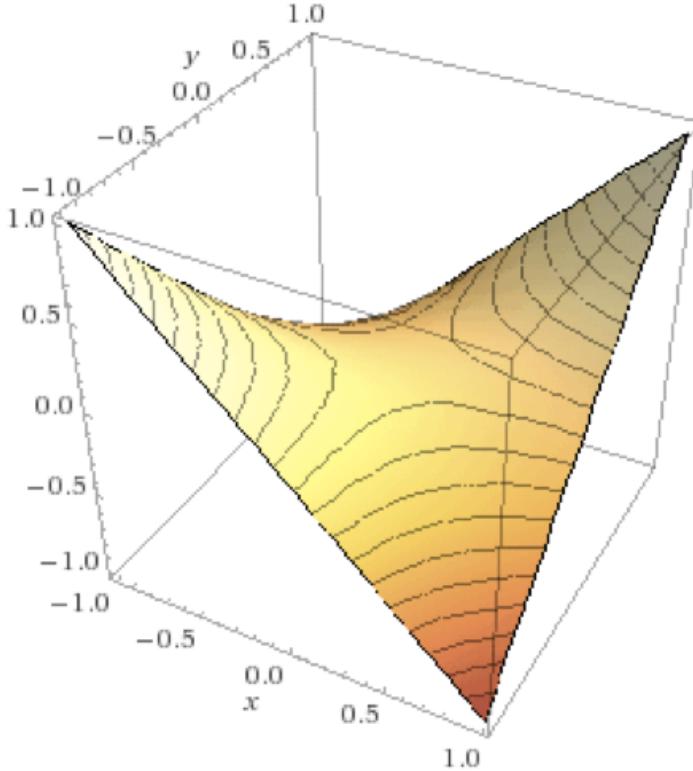


Figure 36: A three-dimensional visualization of the value function $V(x, y) = xy$. This is the canonical example of a function with a saddle point, at $x = y = 0$.

is the simplest possible example of a continuous function with a saddle point. It is easiest to understand this game by visualizing the value function in three dimensions, as shown in figure 36.

The three dimensional visualization shows us clearly that there is a saddle point at $x = y = 0$. This is an equilibrium of the game. We could also have found this point by solving for where the derivatives are zero.

Not every saddle point is an equilibrium; we require that an infinitesimal perturbation of one player's parameters cannot reduce that player's cost. The saddle point for this game satisfies that requirement. It is something of a pathological equilibrium because the value function is constant as a function of each player's parameter when holding the other player's parameter fixed.

To solve for the trajectory taken by gradient descent, we take the derivatives,

and find that

$$\frac{\partial x}{\partial t} = -y(t) \quad (27)$$

$$\frac{\partial y}{\partial t} = x(t). \quad (28)$$

Differentiating equation 28, we obtain

$$\frac{\partial^2 y}{\partial t^2} = \frac{\partial x}{\partial t} = -y(t). \quad (29)$$

Differential equations of this form have sinusoids as their set of basis functions of solutions. Solving for the coefficients that respect the boundary conditions, we obtain

$$x(t) = x(0) \cos(t) - y(0) \sin(t) \quad (30)$$

$$y(t) = x(0) \sin(t) + y(0) \cos(t). \quad (31)$$

These dynamics form a circular orbit, as shown in figure 37. In other words, simultaneous gradient descent with an infinitesimal learning rate will orbit the equilibrium forever, at the same radius that it was initialized. With a larger learning rate, it is possible for simultaneous gradient descent to spiral outward forever. **Simultaneous gradient descent will never approach the equilibrium.**

For some games, simultaneous gradient descent does converge, and for others, such as the one in this exercise, it does not. For GANs, there is no theoretical prediction as to whether simultaneous gradient descent should converge or not. Settling this theoretical question, and developing algorithms guaranteed to converge, remain important open research problems.

8.3 Maximum likelihood in the GAN framework

We wish to find a function f such that the expected gradient of

$$J^{(G)} = \mathbb{E}_{\mathbf{x} \sim p_g} f(\mathbf{x}) \quad (32)$$

is equal to the expected gradient of $D_{\text{KL}}(p_{\text{data}} \| p_g)$.

First we take the derivative of the KL divergence with respect to a parameter θ :

$$\frac{\partial}{\partial \theta} D_{\text{KL}}(p_{\text{data}} \| p_g) = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \frac{\partial}{\partial \theta} \log p_g(\mathbf{x}). \quad (33)$$

We now want to find the f that will make the derivatives of equation 32 match equation 33. We begin by taking the derivatives of equation 32:

$$\frac{\partial}{\partial \theta} J^{(G)} = \mathbb{E}_{\mathbf{x} \sim p_g} f(x) \frac{\partial}{\partial \theta} \log p_g(\mathbf{x}). \quad (34)$$

To obtain this result, we made two assumptions:

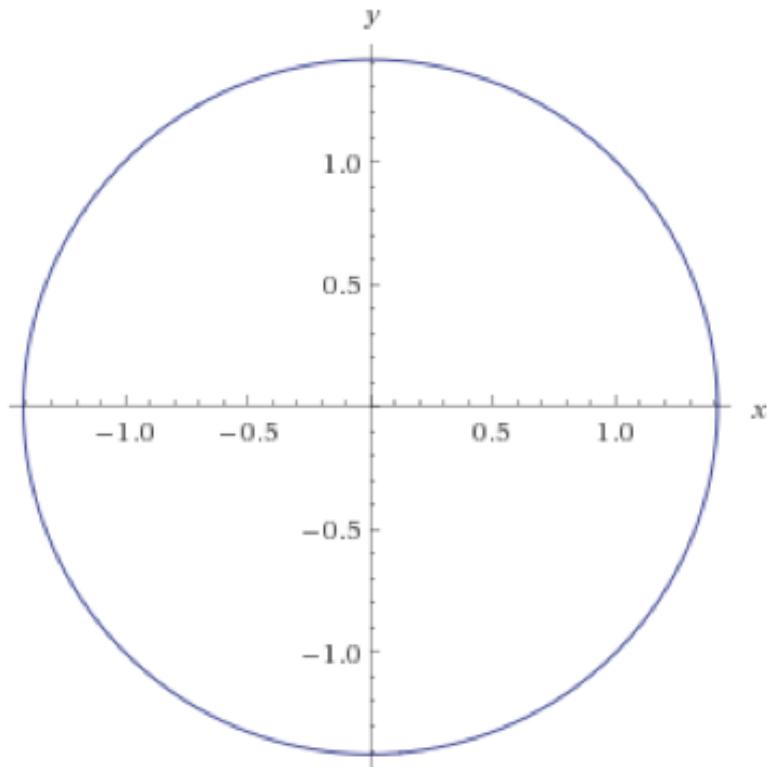


Figure 37: Simultaneous gradient descent with infinitesimal learning rate will orbit indefinitely at constant radius when applied to $V(x, y) = xy$, rather than approaching the equilibrium solution at $x = y = 0$.

1. We assumed that $p_g(\mathbf{x}) \geq 0$ everywhere so that we were able to use the identity $p_g(\mathbf{x}) = \exp(\log p_g(\mathbf{x}))$.
2. We assumed that we can use Leibniz's rule to exchange the order of differentiation and integration (specifically, that both the function and its derivative are continuous, and that the function vanishes for infinite values of \mathbf{x}).

We see that the derivatives of $J^{(G)}$ come very near to giving us what we want; the only problem is that the expectation is computed by drawing samples from p_g when we would like it to be computed by drawing samples from p_{data} . We can fix this problem using an importance sampling trick; by setting $f(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_g(\mathbf{x})}$ we can reweight the contribution to the gradient from each generator sample to compensate for it having been drawn from the generator rather than the data.

Note that when constructing $J^{(G)}$ we must *copy* p_g into $f(\mathbf{x})$ so that $f(\mathbf{x})$ has a derivative of zero with respect to the parameters of p_g . Fortunately, this happens naturally if we obtain the value of $\frac{p_{\text{data}}(\mathbf{x})}{p_g(\mathbf{x})}$.

From section 8.1, we already know that the discriminator estimates the desired ratio. Using some algebra, we can obtain a numerically stable implementation of $f(\mathbf{x})$. If the discriminator is defined to apply a logistic sigmoid function at the output layer, with $D(\mathbf{x}) = \sigma(a(\mathbf{x}))$, then $f(\mathbf{x}) = -\exp(a(\mathbf{x}))$.

This exercise is taken from a result shown by Goodfellow (2014). From this exercise, we see that the discriminator estimates a ratio of densities that can be used to calculate a variety of divergences.

9 Conclusion

GANs are generative models that use supervised learning to approximate an intractable cost function, much as Boltzmann machines use Markov chains to approximate their cost and VAEs use the variational lower bound to approximate their cost. GANs can use this supervised ratio estimation technique to approximate many cost functions, including the KL divergence used for maximum likelihood estimation.

GANs are relatively new and still require some research to reach their new potential. In particular, training GANs requires finding Nash equilibria in high-dimensional, continuous, non-convex games. Researchers should strive to develop better theoretical understanding and better training algorithms for this scenario. Success on this front would improve many other applications, besides GANs.

GANs are crucial to many different state of the art image generation and manipulation systems, and have the potential to enable many other applications in the future.

Acknowledgments

The author would like to thank the NIPS organizers for inviting him to present this tutorial. Many thanks also to those who commented on his Twitter and Facebook posts asking which topics would be of interest to the tutorial audience. Thanks also to D. Kingma for helpful discussions regarding the description of VAEs. Thanks to Zhu Xiaohu, Alex Kurakin and Ilya Edrenkin for spotting typographical errors in the manuscript.

References

- Abadi, M. and Andersen, D. G. (2016). Learning to protect communications with adversarial neural cryptography. *arXiv preprint arXiv:1610.06918*.
- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.
- Ackley, D. H., Hinton, G. E., and Sejnowski, T. J. (1985). A learning algorithm for Boltzmann machines. *Cognitive Science*, **9**, 147–169.
- Bengio, Y., Thibodeau-Laufer, E., Alain, G., and Yosinski, J. (2014). Deep generative stochastic networks trainable by backprop. In *ICML '2014*.
- Brock, A., Lim, T., Ritchie, J. M., and Weston, N. (2016). Neural photo editing with introspective adversarial networks. *CoRR*, **abs/1609.07093**.
- Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., and Abbeel, P. (2016a). Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2172–2180.
- Chen, X., Kingma, D. P., Salimans, T., Duan, Y., Dhariwal, P., Schulman, J., Sutskever, I., and Abbeel, P. (2016b). Variational lossy autoencoder. *arXiv preprint arXiv:1611.02731*.
- Deco, G. and Brauer, W. (1995). Higher order statistical decorrelation without information loss. *NIPS*.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*.

- Deng, J., Berg, A. C., Li, K., and Fei-Fei, L. (2010). What does classifying more than 10,000 image categories tell us? In *Proceedings of the 11th European Conference on Computer Vision: Part V*, ECCV’10, pages 71–84, Berlin, Heidelberg. Springer-Verlag.
- Denton, E., Chintala, S., Szlam, A., and Fergus, R. (2015). Deep generative image models using a Laplacian pyramid of adversarial networks. *NIPS*.
- Dinh, L., Krueger, D., and Bengio, Y. (2014). NICE: Non-linear independent components estimation. arXiv:1410.8516.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. (2016). Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*.
- Donahue, J., Krähenbühl, P., and Darrell, T. (2016). Adversarial feature learning. *arXiv preprint arXiv:1605.09782*.
- Dumoulin, V., Belghazi, I., Poole, B., Lamb, A., Arjovsky, M., Mastropietro, O., and Courville, A. (2016). Adversarially learned inference. *arXiv preprint arXiv:1606.00704*.
- Dziugaite, G. K., Roy, D. M., and Ghahramani, Z. (2015). Training generative neural networks via maximum mean discrepancy optimization. *arXiv preprint arXiv:1505.03906*.
- Edwards, H. and Storkey, A. (2015). Censoring representations with an adversary. *arXiv preprint arXiv:1511.05897*.
- Fahlman, S. E., Hinton, G. E., and Sejnowski, T. J. (1983). Massively parallel architectures for AI: NETL, thistle, and Boltzmann machines. In *Proceedings of the National Conference on Artificial Intelligence AAAI-83*.
- Finn, C. and Levine, S. (2016). Deep visual foresight for planning robot motion. *arXiv preprint arXiv:1610.00696*.
- Finn, C., Christiano, P., Abbeel, P., and Levine, S. (2016a). A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *arXiv preprint arXiv:1611.03852*.
- Finn, C., Goodfellow, I., and Levine, S. (2016b). Unsupervised learning for physical interaction through video prediction. *NIPS*.
- Frey, B. J. (1998). *Graphical models for machine learning and digital communication*. MIT Press.
- Frey, B. J., Hinton, G. E., and Dayan, P. (1996). Does the wake-sleep algorithm learn good density estimators? In D. Touretzky, M. Mozer, and M. Hasselmo, editors, *Advances in Neural Information Processing Systems 8 (NIPS’95)*, pages 661–670. MIT Press, Cambridge, MA.

- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. (2015). Domain-adversarial training of neural networks. *arXiv preprint arXiv:1505.07818*.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- Goodfellow, I. J. (2014). On distinguishability criteria for estimating generative models. In *International Conference on Learning Representations, Workshops Track*.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. (2014a). Explaining and harnessing adversarial examples. *CoRR*, **abs/1412.6572**.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014b). Generative adversarial networks. In *NIPS'2014*.
- Gutmann, M. and Hyvärinen, A. (2010). Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of The Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS'10)*.
- Hinton, G. E. (2007). Learning multiple layers of representation. *Trends in cognitive sciences*, **11**(10), 428–434.
- Hinton, G. E. and Sejnowski, T. J. (1986). Learning and relearning in Boltzmann machines. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, chapter 7, pages 282–317. MIT Press, Cambridge.
- Hinton, G. E., Sejnowski, T. J., and Ackley, D. H. (1984). Boltzmann machines: Constraint satisfaction networks that learn. Technical Report TR-CMU-CS-84-119, Carnegie-Mellon University, Dept. of Computer Science.
- Hinton, G. E., Osindero, S., and Teh, Y. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, **18**, 1527–1554.
- Ho, J. and Ermon, S. (2016). Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*, pages 4565–4573.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift.
- Isola, P., Zhu, J.-Y., Zhou, T., and Efros, A. A. (2016). Image-to-image translation with conditional adversarial networks. *arXiv preprint arXiv:1611.07004*.
- Jang, E., Gu, S., and Poole, B. (2016). Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*.

- Kingma, D. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kingma, D. P. (2013). Fast gradient-based inference with continuous latent variable models in auxiliary form. Technical report, arxiv:1306.0733.
- Kingma, D. P., Salimans, T., and Welling, M. (2016). Improving variational inference with inverse autoregressive flow. *NIPS*.
- Ledig, C., Theis, L., Huszar, F., Caballero, J., Aitken, A. P., Tejani, A., Totz, J., Wang, Z., and Shi, W. (2016). Photo-realistic single image super-resolution using a generative adversarial network. *CoRR*, **abs/1609.04802**.
- Li, Y., Swersky, K., and Zemel, R. S. (2015). Generative moment matching networks. *CoRR*, **abs/1502.02761**.
- Lotter, W., Kreiman, G., and Cox, D. (2015). Unsupervised learning of visual structure using predictive generative networks. *arXiv preprint arXiv:1511.06380*.
- Maddison, C. J., Mnih, A., and Teh, Y. W. (2016). The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*.
- Metz, L., Poole, B., Pfau, D., and Sohl-Dickstein, J. (2016). Unrolled generative adversarial networks. *arXiv preprint arXiv:1611.02163*.
- Nguyen, A., Yosinski, J., Bengio, Y., Dosovitskiy, A., and Clune, J. (2016). Plug & play generative networks: Conditional iterative generation of images in latent space. *arXiv preprint arXiv:1612.00005*.
- Nowozin, S., Cseke, B., and Tomioka, R. (2016). f-gan: Training generative neural samplers using variational divergence minimization. *arXiv preprint arXiv:1606.00709*.
- Odena, A. (2016). Semi-supervised learning with generative adversarial networks. *arXiv preprint arXiv:1606.01583*.
- Oord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. (2016). Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*.
- Pfau, D. and Vinyals, O. (2016). Connecting generative adversarial networks and actor-critic methods. *arXiv preprint arXiv:1610.01945*.
- Radford, A., Metz, L., and Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.

- Ratliff, L. J., Burden, S. A., and Sastry, S. S. (2013). Characterization and computation of local nash equilibria in continuous games. In *Communication, Control, and Computing (Allerton), 2013 51st Annual Allerton Conference on*, pages 917–924. IEEE.
- Reed, S., van den Oord, A., Kalchbrenner, N., Bapst, V., Botvinick, M., and de Freitas, N. (2016a). Generating interpretable images with controllable structure. Technical report.
- Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., and Lee, H. (2016b). Generative adversarial text to image synthesis. *arXiv preprint arXiv:1605.05396*.
- Rezende, D. J. and Mohamed, S. (2015). Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*.
- Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *ICML’2014*. Preprint: arXiv:1401.4082.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. (2014). ImageNet Large Scale Visual Recognition Challenge.
- Salakhutdinov, R. and Hinton, G. (2009). Deep Boltzmann machines. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 5, pages 448–455.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. (2016). Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2226–2234.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of go with deep neural networks and tree search. *Nature*, **529**(7587), 484–489.
- Springenberg, J. T. (2015). Unsupervised and semi-supervised learning with categorical generative adversarial networks. *arXiv preprint arXiv:1511.06390*.
- Springenberg, J. T., Dosovitskiy, A., Brox, T., and Riedmiller, M. (2015). Striving for simplicity: The all convolutional net. In *ICLR*.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. J., and Fergus, R. (2014). Intriguing properties of neural networks. *ICLR, abs/1312.6199*.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2015). Rethinking the Inception Architecture for Computer Vision. *ArXiv e-prints*.

- Theis, L., van den Oord, A., and Bethge, M. (2015). A note on the evaluation of generative models. *arXiv:1511.01844*.
- Warde-Farley, D. and Goodfellow, I. (2016). Adversarial perturbations of deep neural networks. In T. Hazan, G. Papandreou, and D. Tarlow, editors, *Perturbations, Optimization, and Statistics*, chapter 11. MIT Press.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms connectionist reinforcement learning. *Machine Learning*, **8**, 229–256.
- Wu, Y., Burda, Y., Salakhutdinov, R., and Grosse, R. (2016). On the quantitative analysis of decoder-based generative models. *arXiv preprint arXiv:1611.04273*.
- Zhang, H., Xu, T., Li, H., Zhang, S., Huang, X., Wang, X., and Metaxas, D. (2016). Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. *arXiv preprint arXiv:1612.03242*.
- Zhu, J.-Y., Krähenbühl, P., Shechtman, E., and Efros, A. A. (2016). Generative visual manipulation on the natural image manifold. In *European Conference on Computer Vision*, pages 597–613. Springer.