

# Generative Adversarial Nets

✂ Early Version In Arxiv, new version in NeurIPS

Ian J. Goodfellow, Jean Pouget-Abadie\*, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair†, Aaron Courville, Yoshua Bengio‡

Département d'informatique et de recherche opérationnelle  
Université de Montréal  
Montréal, QC H3C 3J7

- Unsupervised learning
- Use supervised learning loss function to do unsupervised learning (supervised learning loss function make train process more efficient) (inspired models like BERT)

## Abstract

✂ 主要关注GAN在于什么,而不在于与别人工作的区别 ← 开创性论文写法

We propose a new framework for estimating generative models via an **adversarial process**, in which we **simultaneously train two models: a generative model  $G$  that captures the data distribution, and a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than  $G$ . The training procedure for  $G$  is to maximize the probability of  $D$  making a mistake.** This framework corresponds to a **minimax two-player game**. In the space of arbitrary functions  $G$  and  $D$ , **a unique solution exists, with  $G$  recovering the training data distribution and  $D$  equal to  $\frac{1}{2}$  everywhere.** In the case where  $G$  and  $D$  are defined by **multilayer perceptrons**, the entire system can be trained with **backpropagation**. There is no need for any Markov chains or unrolled approximate inference networks during either training or generation of samples. Experiments demonstrate the potential of the framework through qualitative and quantitative evaluation of the generated samples.

model data distribution to generate data distribution s.t. the data sampled from the new distribution has nearly no difference from the real data (i.e. training data), which means making discriminative model  $D$  hard, even not able to check whether sample is from read data

function space

training data distribution = true distribution

误差的反向传递

## 1 Introduction

Objective

The promise of deep learning is to discover rich, hierarchical models [2] that represent probability distributions over the kinds of data encountered in artificial intelligence applications, such as natural images, audio waveforms containing speech, and symbols in natural language corpora. So far, the most striking successes in deep learning have involved discriminative models, usually those that map a high-dimensional, rich sensory input to a class label [14, 22]. These striking successes have primarily been based on the backpropagation and dropout algorithms, using piecewise linear units [19, 9, 10] which have a particularly well-behaved gradient. Deep *generative* models have had less of an impact, due to **the difficulty of approximating many intractable probabilistic computations that arise in maximum likelihood estimation and related strategies, and due to difficulty of leveraging the benefits of piecewise linear units in the generative context.** We propose a new generative model estimation procedure that sidesteps these difficulties.<sup>1</sup>

deep learning is not just about deep neural nets, but a representation of data distribution features

深度学习在生成模型上进展不多

↑  
在最大化似然函数时需要将概率分布进行很多近似,近似问题带来了许多计算上的困难。

↓  
不是通过近似似然函数而是通过其它方法来获得计算上更好的模型。

In the proposed *adversarial nets* framework, the generative model is pitted against an adversary: a discriminative model that learns to determine whether a **sample is from the model distribution or the data distribution.** The generative model can be thought of as analogous to a team of counterfeiters, trying to produce fake currency and use it without detection, while the discriminative model is analogous to the police, trying to detect the counterfeit currency. Competition in this game drives both teams to improve their methods until the counterfeits are indistinguishable from the genuine articles.

Analogy

\*Jean Pouget-Abadie is visiting Université de Montréal from Ecole Polytechnique.

†Sherjil Ozair is visiting Université de Montréal from Indian Institute of Technology Delhi

‡Yoshua Bengio is a CIFAR Senior Fellow.

<sup>1</sup>All code and hyperparameters available at <http://www.github.com/goodfeli/adversarial>

对抗网络定义

This framework can yield specific training algorithms for many kinds of model and optimization algorithm. In this article, we explore the special case when the generative model generates samples by passing random noise through a multilayer perceptron, and the discriminative model is also a multilayer perceptron. We refer to this special case as *adversarial nets*. In this case, we can train both models using only the highly successful backpropagation and dropout algorithms [17] and sample from the generative model using only forward propagation. No approximate inference or Markov chains are necessary.

输入随机噪声,MLP能够将产生随机噪声的分布(Commonly to be Gaussian)映射到任意一个我们想拟合的分布

因为两个模型都是基于MLP,所以在训练时可以直接通过误差的backprop,从而不需要通过类似Markov chain一样的复杂算法来对分布进行复杂的采样,使GAN在计算上有优势

## 2 Related work

An alternative to directed graphical models with latent variables are undirected graphical models with latent variables, such as restricted Boltzmann machines (RBMs) [27, 16], deep Boltzmann machines (DBMs) [26] and their numerous variants. The interactions within such models are represented as the product of unnormalized potential functions, normalized by a global summation/integration over all states of the random variables. This quantity (the *partition function*) and its gradient are intractable for all but the most trivial instances, although they can be estimated by Markov chain Monte Carlo (MCMC) methods. Mixing poses a significant problem for learning algorithms that rely on MCMC [3, 5].

Deep belief networks (DBNs) [16] are hybrid models containing a single undirected layer and several directed layers. While a fast approximate layer-wise training criterion exists, DBNs incur the computational difficulties associated with both undirected and directed models.

Alternative criteria that do not approximate or bound the log-likelihood have also been proposed, such as score matching [18] and noise-contrastive estimation (NCE) [13]. Both of these require the learned probability density to be analytically specified up to a normalization constant. Note that in many interesting generative models with several layers of latent variables (such as DBNs and DBMs), it is not even possible to derive a tractable unnormalized probability density. Some models such as denoising auto-encoders [30] and contractive autoencoders have learning rules very similar to score matching applied to RBMs. In NCE, as in this work, a discriminative training criterion is employed to fit a generative model. However, rather than fitting a separate discriminative model, the generative model itself is used to discriminate generated data from samples a fixed noise distribution. Because NCE uses a fixed noise distribution, learning slows dramatically after the model has learned even an approximately correct distribution over a small subset of the observed variables.

Finally, some techniques do not involve defining a probability distribution explicitly, but rather train a generative machine to draw samples from the desired distribution. This approach has the advantage that such machines can be designed to be trained by back-propagation. Prominent recent work in this area includes the generative stochastic network (GSN) framework [5], which extends generalized denoising auto-encoders [4]: both can be seen as defining a parameterized Markov chain, i.e., one learns the parameters of a machine that performs one step of a generative Markov chain. Compared to GSNs, the adversarial nets framework does not require a Markov chain for sampling. Because adversarial nets do not require feedback loops during generation, they are better able to leverage piecewise linear units [19, 9, 10], which improve the performance of backpropagation but have problems with unbounded activation when used in a feedback loop. More recent examples of training a generative machine by back-propagating into it include recent work on auto-encoding variational Bayes [20] and stochastic backpropagation [24].

## 3 Adversarial nets (Model, Objective function & Optimization)

The adversarial modeling framework is most straightforward to apply when the models are both multilayer perceptrons. To learn the generator's distribution  $p_g$  over data  $x$ , we define a prior on input noise variables  $p_z(z)$ , then represent a mapping to data space as  $G(z; \theta_g)$ , where  $G$  is a differentiable function represented by a multilayer perceptron with parameters  $\theta_g$ . We also define a second multilayer perceptron  $D(x; \theta_d)$  that outputs a single scalar.  $D(x)$  represents the probability that  $x$  came from the data rather than  $p_g$ . We train  $D$  to maximize the probability of assigning the correct label to both training examples and samples from  $G$ . We simultaneously train  $G$  to minimize  $\log(1 - D(G(z)))$ :

$\theta_g$  为可学习参数

$\theta_d$  为可学习参数

the output scalar used to discriminate whether  $x$  from real data or from  $G$

if  $D$  is very good (i.e.) correctly discriminate every  $x$  (i.e.  $D$  believes  $x$  from  $G$ ), then  $D(G(x)) \rightarrow 0$  such that  $\log(1 - D(G)) \rightarrow \log(1 - 0) = \log(1) = 0$ , otherwise it will output number  $> 0$ , to an extreme case  $D(G(x)) \rightarrow 1$ , which means  $D$  believes  $x$  from real data, s.t.  $\log(1 - D(G)) \rightarrow \log(1 - 1) = \log(0) = -\infty$

★ (a) is expectation where  $x$  sampled from real data, assume  $D$  is perfect then,  $D(x) \rightarrow 1$ ,  $\log(D(x)) \rightarrow 0$ ; (b) is expectation where  $x$  sampled from noise distribution  $p_z(z)$ , put into generator  $G$  to generate  $x$ . If  $D$  is perfect then  $D(G(z)) \rightarrow 0$ , and  $\log(1-0)=0$ ; otherwise if  $D$  is not perfect, then both (a) and (b) can be negative, in order to make best  $D \in [0,1]$ , we want to maximize  $V$ . Since we want  $D$  unable to discriminate  $x$  from  $G$ , we need train  $G$  s.t.  $D(G(z)) \rightarrow 1$ , which means  $D$  making mistakes, then  $\log(1-D(G(z))) \rightarrow -\infty$ , meaning  $V$  needs to be minimized by  $G$ .

in game theory, as we reach a point (equilibrium, or Nash Equilibrium) where  $D$  and  $G$  can not improve anymore.

In other words,  $D$  and  $G$  play the following two-player minimax game with value function  $V(G, D)$ :

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (1)$$

In the next section, we present a theoretical analysis of adversarial nets, essentially showing that the training criterion allows one to recover the data generating distribution as  $G$  and  $D$  are given enough capacity, i.e., in the non-parametric limit. See Figure 1 for a less formal, more pedagogical explanation of the approach. In practice, we must implement the game using an iterative, numerical approach. Optimizing  $D$  to completion in the inner loop of training is computationally prohibitive, and on finite datasets would result in overfitting. Instead, we alternate between  $k$  steps of optimizing  $D$  and one step of optimizing  $G$ . This results in  $D$  being maintained near its optimal solution, so long as  $G$  changes slowly enough. This strategy is analogous to the way that SML/PCD [31, 29] training maintains samples from a Markov chain from one learning step to the next in order to avoid burning in a Markov chain as part of the inner loop of learning. The procedure is formally presented in Algorithm 1.

In practice, equation 1 may not provide sufficient gradient for  $G$  to learn well. Early in learning, when  $G$  is poor,  $D$  can reject samples with high confidence because they are clearly different from the training data. In this case,  $\log(1 - D(G(z)))$  saturates. Rather than training  $G$  to minimize  $\log(1 - D(G(z)))$  we can train  $G$  to maximize  $\log D(G(z))$ . This objective function results in the same fixed point of the dynamics of  $G$  and  $D$  but provides much stronger gradients early in learning.

★ 第二项(b)会有问题,早期 $G$ 较弱,生成的数据与真实的数据相差很远,使 $D$ 很容易训练得特别好,产生(下项笔记中提到的)问题,所以这里作者建议将目标函数改成最大化这一项。

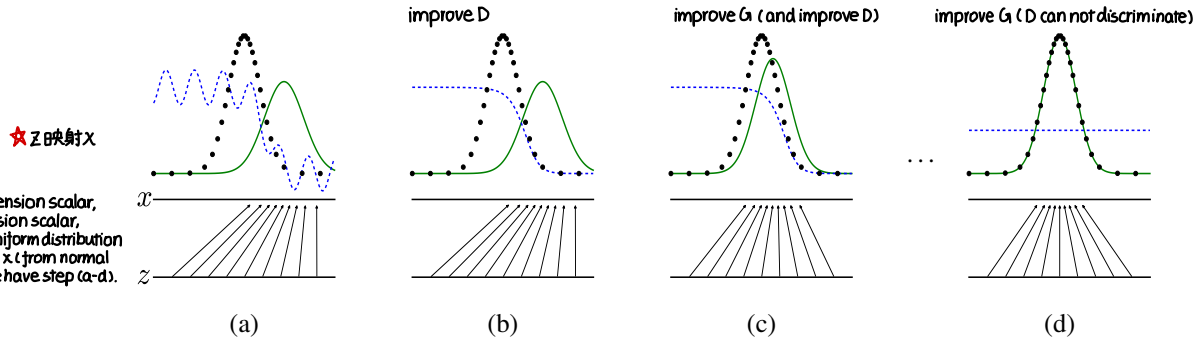


Figure 1: Generative adversarial nets are trained by simultaneously updating the discriminative distribution ( $D$ , blue, dashed line) so that it discriminates between samples from the data generating distribution (black, dotted line)  $p_x$  from those of the generative distribution  $p_g$  ( $G$ ) (green, solid line). The lower horizontal line is the domain from which  $z$  is sampled, in this case uniformly. The horizontal line above is part of the domain of  $x$ . The upward arrows show how the mapping  $x = G(z)$  imposes the non-uniform distribution  $p_g$  on transformed samples.  $G$  contracts in regions of high density and expands in regions of low density of  $p_g$ . (a) Consider an adversarial pair near convergence:  $p_g$  is similar to  $p_{\text{data}}$  and  $D$  is a partially accurate classifier. (b) In the inner loop of the algorithm  $D$  is trained to discriminate samples from data, converging to  $D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$ . (c) After an update to  $G$ , gradient of  $D$  has guided  $G(z)$  to flow to regions that are more likely to be classified as data. (d) After several steps of training, if  $G$  and  $D$  have enough capacity, they will reach a point at which both cannot improve because  $p_g = p_{\text{data}}$ . The discriminator is unable to differentiate between the two distributions, i.e.  $D(x) = \frac{1}{2}$ .

## 4 Theoretical Results

The generator  $G$  implicitly defines a probability distribution  $p_g$  as the distribution of the samples  $G(z)$  obtained when  $z \sim p_z$ . Therefore, we would like Algorithm 1 to converge to a good estimator of  $p_{\text{data}}$ , if given enough capacity and training time. The results of this section are done in a non-parametric setting, e.g. we represent a model with infinite capacity by studying convergence in the space of probability density functions.

We will show in section 4.1 that this minimax game has a global optimum for  $p_g = p_{\text{data}}$ . We will then show in section 4.2 that Algorithm 1 optimizes Eq 1, thus obtaining the desired result.

理论结论的两个构成部分

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

**for** number of training iterations **do** ← 迭代次数直到完成 How to know the alg converge? ←

★ **for**  $k$  steps **do**  $k$  is hyper-parameter,  $k$  不能太小也不能太大 → 保证  $D$  有足够的更新, 但也不要更新得太好, 使得  $D$  更新与  $G$  更新在进度上差不多

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{\text{data}}(x)$ .
- Update the discriminator by ascending its stochastic gradient:  
放进价值函数中求梯度, 对  $D$  求导数来更新  $D$ , 做  $k$  步

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right]$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Update the generator by descending its stochastic gradient:  
生成器  $G$  与第一项无关, 所以只计算第二项, 更新生成器

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)})))$$

★ 若  $D$  未更新好, 但  $G$  已改变, 再用未更新的  $D$  去辨别  $G$  则意义不大; 反之更新就把  $D$  训练到完美, 则  $D(G(z)) \rightarrow 0, \log(1-0)=0$ , 对  $D$  求导 (求梯度) 会在你的生成模型上更新产生困难 (梯度消失), 直觉上,  $D$  过强会使  $G$  找不到学习改进的方向,  $D$  太弱会使  $G$  改进过早收尾。

**end for**

Very unstable convergence  
若一方不动, 另一方变乱算不算收敛。  
若两方来回抖动算不算收敛。

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

#### 4.1 Global Optimality of $p_g = p_{\text{data}}$

We first consider the optimal discriminator  $D$  for any given generator  $G$ .

**Proposition 1.** For  $G$  fixed, the optimal discriminator  $D$  is

① Optimize/Solve  $\max_D$

$$\text{最优解 } D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \in [0, 1] \quad (2)$$

when  $p_{\text{data}} = p_g$ , we have  $D_G^*(x) = \frac{1}{2}$  表示完全无法区分这两个分布

★ Two-sample Test, 判断两块数据是否来自同一分布 (统计中还有其他工具用于解决这一问题, 如下 Test), 这里可以完全不管分布是什么样子 (无视在高维上很多统计工具不好用), 直接训练一个二分类器, 若它能分开这两块数据, 就表示它们来自不同分布, 反之则同分布。

*Proof.* The training criterion for the discriminator  $D$ , given any generator  $G$ , is to maximize the quantity  $V(G, D)$

※ Calculate Expectation:  $\mathbb{E}_{x \sim p}(x) = \int x p(x) dx$

$$\begin{aligned} V(G, D) &= \int_x p_{\text{data}}(x) \log(D(x)) dx + \int_z p_g(z) \log(1 - D(G(z))) dz \\ &= \int_x p_{\text{data}}(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx \end{aligned}$$

$\mathbb{E}_{x \sim p_{\text{data}}}[a(x)]$   $\mathbb{E}_{z \sim p_g}[b(z)]$   
 $x \rightarrow x$  的分布来源于生成器所对应的  $p_g$ , 替换  $z$  后两项合并

$$a = p_{\text{data}}(x) \in [0, 1], b = p_g(x) \in [0, 1], y = D(x), \text{ result in a convex function about } y$$

$$\frac{dy}{dy} = \frac{a}{y} - \frac{b}{1-y} = 0 \Rightarrow y = \frac{a}{a+b}$$

For any  $(a, b) \in \mathbb{R}^2 \setminus \{0, 0\}$ , the function  $y \rightarrow a \log(y) + b \log(1 - y)$  achieves its maximum in  $[0, 1]$  at  $\frac{a}{a+b}$ . The discriminator does not need to be defined outside of  $\text{Supp}(p_{\text{data}}) \cup \text{Supp}(p_g)$ , concluding the proof.  $\square$

this proves equation (2)

Note that the training objective for  $D$  can be interpreted as maximizing the log-likelihood for estimating the conditional probability  $P(Y = y|x)$ , where  $Y$  indicates whether  $x$  comes from  $p_{\text{data}}$  (with  $y = 1$ ) or from  $p_g$  (with  $y = 0$ ). The minimax game in Eq. 1 can now be reformulated as:

将  $D$  最优解代入  $V(G, D)$  可得关于  $G$  的函数

$$\begin{aligned} C(G) &= \max_D V(G, D) \\ &= \mathbb{E}_{x \sim p_{\text{data}}} [\log D_G^*(x)] + \mathbb{E}_{z \sim p_g} [\log(1 - D_G^*(G(z)))] \\ &= \mathbb{E}_{x \sim p_{\text{data}}} [\log D_G^*(x)] + \mathbb{E}_{x \sim p_g} [\log(1 - D_G^*(x))] \\ &= \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right] + \mathbb{E}_{x \sim p_g} \left[ \log \frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right] \end{aligned} \quad (4)$$

※信息论概念: KL散度用于衡量两个分布的相似度。  
 $KL(p||q) = \mathbb{E}_{x \sim p} \log(p(x)/q(x))$ , 表示在已知p分布的情况下, 至少多少比特能把q分布描述出来

★  $p_{data} + p_g$  不是一个分布, 因为它求积后等于2, 但是  $\frac{1}{2}(p_{data}(x) + p_g(x))$  积为1可构成一个分布, 也就是:  
 $\log(p_{data}(x)/[\frac{1}{2}(p_{data}(x) + p_g(x))]) + \log(2) - \log(2)$   
 $= \log(2 \cdot p_{data}(x)/[p_{data}(x) + p_g(x)]) - \log(2)$   
 $= \log(p_{data}(x)/[\frac{1}{2}(p_{data}(x) + p_g(x))]) - \log(2)$   
 由此上下能构成两个分布, 第一个log项可构成KL散度。  
 $= KL(p_{data}(x) || [\frac{1}{2}(p_{data}(x) + p_g(x))]) - \log(2)$

② Optimize/Solve min Theorem 1. The global minimum of the virtual training criterion  $C(G)$  is achieved if and only if  $p_g = p_{data}$ . At that point,  $C(G)$  achieves the value  $-\log 4$ .

Proof. For  $p_g = p_{data}$ ,  $D_G^*(x) = \frac{1}{2}$ , (consider Eq. 2). Hence, by inspecting Eq. 4 at  $D_G^*(x) = \frac{1}{2}$ , we find  $C(G) = \log \frac{1}{2} + \log \frac{1}{2} = -\log 4$ . To see that this is the best possible value of  $C(G)$ , reached only for  $p_g = p_{data}$ , observe that

$$\mathbb{E}_{x \sim p_{data}} [-\log 2] + \mathbb{E}_{x \sim p_g} [-\log 2] = -\log 4$$

and that by subtracting this expression from  $C(G) = V(D_G^*, G)$ , we obtain:

$$C(G) = -\log(4) + KL\left(p_{data} \left\| \frac{p_{data} + p_g}{2}\right.\right) + KL\left(p_g \left\| \frac{p_{data} + p_g}{2}\right.\right) \quad (5)$$

★ KL散度一定 $\geq 0$ , 当 $KL(p||q) = 0$ 时,  $p = q$  也就是说“当” $C(G)$ 取最小要使两个KL散度项为0, 意味着  $p_{data} = (p_{data} + p_g)/2$ , 即  $p_{data} = p_g$

where KL is the Kullback-Leibler divergence. We recognize in the previous expression the Jensen-Shannon divergence between the model's distribution and the data generating process:

$$C(G) = -\log(4) + 2 \cdot JSD(p_{data} || p_g) \quad (6)$$

※ JSD相对KLD, 不同在它是对称的(KL不能使p与q互换而得到相同结果, 而JSD可以←有评价标准因为GAN使用JSD, 使得计算更简单。

Since the Jensen-Shannon divergence between two distributions is always non-negative and zero only when they are equal, we have shown that  $C^* = -\log(4)$  is the global minimum of  $C(G)$  and that the only solution is  $p_g = p_{data}$ , i.e., the generative model perfectly replicating the data generating process. □

## 4.2 Convergence of Algorithm 1 Optimization of alg 1

Proposition 2. If  $G$  and  $D$  have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given  $G$ , and  $p_g$  is updated so as to improve the criterion

每一步D可以达到它的最优解

$$\mathbb{E}_{x \sim p_{data}} [\log D_G^*(x)] + \mathbb{E}_{x \sim p_g} [\log(1 - D_G^*(x))]$$

对G的优化是去迭代这个步骤。

then  $p_g$  converges to  $p_{data}$

G已换成最优解

★ 将价值函数看成一个关于 $p_g$ 的函数, ( $p_g$  is our model, or distribution, or function) 即关于函数的函数。  $U(p_g, D)$  is convex in  $p_g$ . 一个凸函数的上限函数还是凸函数  $\sup_{\alpha \in \mathcal{A}} U(p_g, D)$  is convex in  $p_g$  with a unique global optima as proven, 即对函数作梯度下降可得最优解

Proof. Consider  $V(G, D) = U(p_g, D)$  as a function of  $p_g$  as done in the above criterion. Note that  $U(p_g, D)$  is convex in  $p_g$ . The subderivatives of a supremum of convex functions include the derivative of the function at the point where the maximum is attained. In other words, if  $f(x) = \sup_{\alpha \in \mathcal{A}} f_{\alpha}(x)$  and  $f_{\alpha}(x)$  is convex in  $x$  for every  $\alpha$ , then  $\partial f_{\beta}(x) \in \partial f$  if  $\beta = \arg \sup_{\alpha \in \mathcal{A}} f_{\alpha}(x)$ . This is equivalent to computing a gradient descent update for  $p_g$  at the optimal  $D$  given the corresponding  $G$ .  $\sup_D U(p_g, D)$  is convex in  $p_g$  with a unique global optima as proven in Thm 1, therefore with sufficiently small updates of  $p_g$ ,  $p_g$  converges to  $p_x$ , concluding the proof. □

★ 虽然以上假设了每次迭代都会对D优化到极致, 但在实际算法中并没有, 只是迭代1k步, 所以该证明无法说明算法可行。(实际上很难收敛)

In practice, adversarial nets represent a limited family of  $p_g$  distributions via the function  $G(z; \theta_g)$ , and we optimize  $\theta_g$  rather than  $p_g$  itself. Using a multilayer perceptron to define  $G$  introduces multiple critical points in parameter space. However, the excellent performance of multilayer perceptrons in practice suggests that they are a reasonable model to use despite their lack of theoretical guarantees.

## 5 Experiments

We trained adversarial nets on a range of datasets including MNIST[23], the Toronto Face Database (TFD) [28], and CIFAR-10 [21]. The generator nets used a mixture of rectifier linear activations [19, 9] and sigmoid activations, while the discriminator net used maxout [10] activations. Dropout [17] was applied in training the discriminator net. While our theoretical framework permits the use of dropout and other noise at intermediate layers of the generator, we used noise as the input to only the bottommost layer of the generator network.

We estimate probability of the test set data under  $p_g$  by fitting a Gaussian Parzen window to the samples generated with  $G$  and reporting the log-likelihood under this distribution. The  $\sigma$  parameter



Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 66$
Stacked CAE [3]	$121 \pm 1.6$	<b><math>2110 \pm 50</math></b>
Deep GSN [6]	$214 \pm 1.1$	$1890 \pm 29$
Adversarial nets	<b><math>225 \pm 2</math></b>	<b><math>2057 \pm 26</math></b>

Table 1: Parzen window-based log-likelihood estimates. The reported numbers on MNIST are the mean log-likelihood of samples on test set, with the standard error of the mean computed across examples. On TFD, we computed the standard error across folds of the dataset, with a different  $\sigma$  chosen using the validation set of each fold. On TFD,  $\sigma$  was cross validated on each fold and mean log-likelihood on each fold were computed. For MNIST we compare against other models of the real-valued (rather than binary) version of dataset.

of the Gaussians was obtained by cross validation on the validation set. This procedure was introduced in Breuleux *et al.* [8] and used for various generative models for which the exact likelihood is not tractable [25, 3, 5]. Results are reported in Table 1. This method of estimating the likelihood has somewhat high variance and does not perform well in high dimensional spaces but it is the best method available to our knowledge. Advances in generative models that can sample but not estimate likelihood directly motivate further research into how to evaluate such models.

In Figures 2 and 3 we show samples drawn from the generator net after training. While we make no claim that these samples are better than samples generated by existing methods, we believe that these samples are at least competitive with the better generative models in the literature and highlight the potential of the adversarial framework.

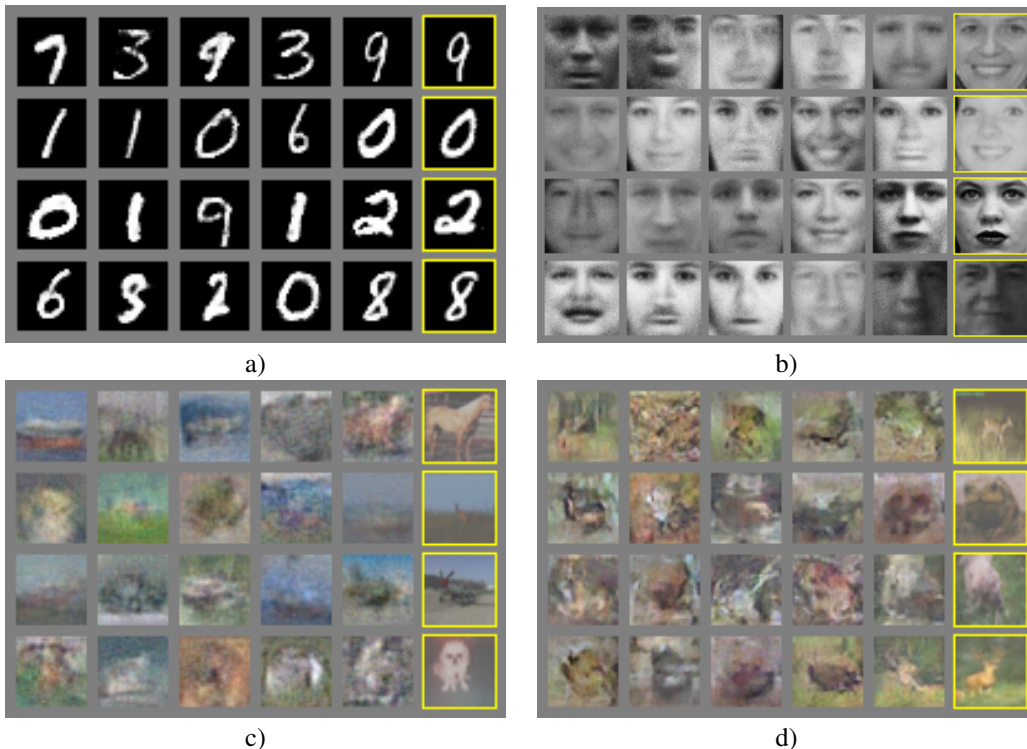


Figure 2: Visualization of samples from the model. Rightmost column shows the nearest training example of the neighboring sample, in order to demonstrate that the model has not memorized the training set. Samples are fair random draws, not cherry-picked. Unlike most other visualizations of deep generative models, these images show actual samples from the model distributions, not conditional means given samples of hidden units. Moreover, these samples are uncorrelated because the sampling process does not depend on Markov chain mixing. a) MNIST b) TFD c) CIFAR-10 (fully connected model) d) CIFAR-10 (convolutional discriminator and “deconvolutional” generator)



Figure 3: Digits obtained by linearly interpolating between coordinates in  $z$  space of the full model.

	Deep directed graphical models	Deep undirected graphical models	Generative autoencoders	Adversarial models
Training	Inference needed during training.	Inference needed during training. MCMC needed to approximate partition function gradient.	Enforced tradeoff between mixing and power of reconstruction generation	Synchronizing the discriminator with the generator. Helvetica.
Inference	Learned approximate inference	Variational inference	MCMC-based inference	Learned approximate inference
Sampling	No difficulties	Requires Markov chain	Requires Markov chain	No difficulties
Evaluating $p(x)$	Intractable, may be approximated with AIS	Intractable, may be approximated with AIS	Not explicitly represented, may be approximated with Parzen density estimation	Not explicitly represented, may be approximated with Parzen density estimation
Model design	Nearly all models incur extreme difficulty	Careful design needed to ensure multiple properties	Any differentiable function is theoretically permitted	Any differentiable function is theoretically permitted

Table 2: Challenges in generative modeling: a summary of the difficulties encountered by different approaches to deep generative modeling for each of the major operations involving a model.

## 6 Advantages and disadvantages

This new framework comes with advantages and disadvantages relative to previous modeling frameworks. The disadvantages are primarily that **there is no explicit representation of  $p_g(x)$ , and that  $D$  must be synchronized well with  $G$  during training** (in particular,  $G$  must not be trained too much without updating  $D$ , in order to avoid “the Helvetica scenario” in which  $G$  collapses too many values of  $z$  to the same value of  $x$  to have enough diversity to model  $p_{\text{data}}$ ), much as the negative chains of a Boltzmann machine must be kept up to date between learning steps. The advantages are that **Markov chains are never needed, only backprop is used to obtain gradients, no inference is needed during learning, and a wide variety of functions can be incorporated into the model**. Table 2 summarizes the comparison of generative adversarial nets with other generative modeling approaches.

**DisADV**  
计算比较困难, G与D必须均衡好  
Model Collapse; Initialization, hard to converge; expensive to adjust hyper-parameters  
**ADV**  
因为生成器并没有去看真实数据, 试图去拟合那些特征, 使得它能生成比较锐利的边缘 (后来被证明)

The aforementioned advantages are primarily computational. Adversarial models may also gain some statistical advantage from the generator network not being updated directly with data examples, but only with gradients flowing through the discriminator. This means that components of the input are not copied directly into the generator’s parameters. Another advantage of adversarial networks is that they can represent very sharp, even degenerate distributions, while methods based on Markov chains require that the distribution be somewhat blurry in order for the chains to be able to mix between modes.

## ✕7 Conclusions and future work

This framework admits many straightforward extensions:

1. A *conditional generative model*  $p(x | c)$  can be obtained by adding  $c$  as input to both  $G$  and  $D$ . 当前不受控制
2. *Learned approximate inference* can be performed by training an auxiliary network to predict  $z$  given  $x$ . This is similar to the inference net trained by the wake-sleep algorithm [15] but with the advantage that the inference net may be trained for a fixed generator net after the generator net has finished training.

3. One can approximately model all conditionals  $p(x_S | x_{\bar{S}})$  where  $S$  is a subset of the indices of  $x$  by training a family of conditional models that share parameters. Essentially, one can use adversarial nets to implement a stochastic extension of the deterministic MP-DBM [11].
4. *Semi-supervised learning*: features from the discriminator or inference net could improve performance of classifiers when limited labeled data is available.
5. *Efficiency improvements*: training could be accelerated greatly by devising better methods for coordinating  $G$  and  $D$  or determining better distributions to sample  $z$  from during training.

This paper has demonstrated the viability of the adversarial modeling framework, suggesting that these research directions could prove useful.

## Acknowledgments

We would like to acknowledge Patrice Marcotte, Olivier Delalleau, Kyunghyun Cho, Guillaume Alain and Jason Yosinski for helpful discussions. Yann Dauphin shared his Parzen window evaluation code with us. We would like to thank the developers of Pylearn2 [12] and Theano [7, 1], particularly Frédéric Bastien who rushed a Theano feature specifically to benefit this project. Arnaud Bergeron provided much-needed support with L<sup>A</sup>T<sub>E</sub>X typesetting. We would also like to thank CIFAR, and Canada Research Chairs for funding, and Compute Canada, and Calcul Québec for providing computational resources. Ian Goodfellow is supported by the 2013 Google Fellowship in Deep Learning. Finally, we would like to thank Les Trois Brasseurs for stimulating our creativity.

## References

- [1] Bastien, F., Lamblin, P., Pascanu, R., Bergstra, J., Goodfellow, I. J., Bergeron, A., Bouchard, N., and Bengio, Y. (2012). Theano: new features and speed improvements. Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop.
- [2] Bengio, Y. (2009). *Learning deep architectures for AI*. Now Publishers.
- [3] Bengio, Y., Mesnil, G., Dauphin, Y., and Rifai, S. (2013a). Better mixing via deep representations. In *ICML'13*.
- [4] Bengio, Y., Yao, L., Alain, G., and Vincent, P. (2013b). Generalized denoising auto-encoders as generative models. In *NIPS26*. Nips Foundation.
- [5] Bengio, Y., Thibodeau-Laufer, E., and Yosinski, J. (2014a). Deep generative stochastic networks trainable by backprop. In *ICML'14*.
- [6] Bengio, Y., Thibodeau-Laufer, E., Alain, G., and Yosinski, J. (2014b). Deep generative stochastic networks trainable by backprop. In *Proceedings of the 30th International Conference on Machine Learning (ICML'14)*.
- [7] Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., Turian, J., Warde-Farley, D., and Bengio, Y. (2010). Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)*. Oral Presentation.
- [8] Breuleux, O., Bengio, Y., and Vincent, P. (2011). Quickly generating representative samples from an RBM-derived process. *Neural Computation*, **23**(8), 2053–2073.
- [9] Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep sparse rectifier neural networks. In *AISTATS'2011*.
- [10] Goodfellow, I. J., Warde-Farley, D., Mirza, M., Courville, A., and Bengio, Y. (2013a). Maxout networks. In *ICML'2013*.
- [11] Goodfellow, I. J., Mirza, M., Courville, A., and Bengio, Y. (2013b). Multi-prediction deep Boltzmann machines. In *NIPS'2013*.
- [12] Goodfellow, I. J., Warde-Farley, D., Lamblin, P., Dumoulin, V., Mirza, M., Pascanu, R., Bergstra, J., Bastien, F., and Bengio, Y. (2013c). Pylearn2: a machine learning research library. *arXiv preprint arXiv:1308.4214*.
- [13] Gutmann, M. and Hyvarinen, A. (2010). Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS'2010*.
- [14] Hinton, G., Deng, L., Dahl, G. E., Mohamed, A., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T., and Kingsbury, B. (2012a). Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal Processing Magazine*, **29**(6), 82–97.
- [15] Hinton, G. E., Dayan, P., Frey, B. J., and Neal, R. M. (1995). The wake-sleep algorithm for unsupervised neural networks. *Science*, **268**, 1558–1561.



- [16] Hinton, G. E., Osindero, S., and Teh, Y. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, **18**, 1527–1554.
- [17] Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2012b). Improving neural networks by preventing co-adaptation of feature detectors. Technical report, arXiv:1207.0580.
- [18] Hyvärinen, A. (2005). Estimation of non-normalized statistical models using score matching. *J. Machine Learning Res.*, **6**.
- [19] Jarrett, K., Kavukcuoglu, K., Ranzato, M., and LeCun, Y. (2009). What is the best multi-stage architecture for object recognition? In *Proc. International Conference on Computer Vision (ICCV'09)*, pages 2146–2153. IEEE.
- [20] Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- [21] Krizhevsky, A. and Hinton, G. (2009). Learning multiple layers of features from tiny images. Technical report, University of Toronto.
- [22] Krizhevsky, A., Sutskever, I., and Hinton, G. (2012). ImageNet classification with deep convolutional neural networks. In *NIPS'2012*.
- [23] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, **86**(11), 2278–2324.
- [24] Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. Technical report, arXiv:1401.4082.
- [25] Rifai, S., Bengio, Y., Dauphin, Y., and Vincent, P. (2012). A generative process for sampling contractive auto-encoders. In *ICML'12*.
- [26] Salakhutdinov, R. and Hinton, G. E. (2009). Deep Boltzmann machines. In *AISTATS'2009*, pages 448–455.
- [27] Smolensky, P. (1986). Information processing in dynamical systems: Foundations of harmony theory. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, chapter 6, pages 194–281. MIT Press, Cambridge.
- [28] Susskind, J., Anderson, A., and Hinton, G. E. (2010). The Toronto face dataset. Technical Report UTML TR 2010-001, U. Toronto.
- [29] Tieleman, T. (2008). Training restricted Boltzmann machines using approximations to the likelihood gradient. In W. W. Cohen, A. McCallum, and S. T. Roweis, editors, *ICML 2008*, pages 1064–1071. ACM.
- [30] Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P.-A. (2008). Extracting and composing robust features with denoising autoencoders. In *ICML 2008*.
- [31] Yunes, L. (1999). On the convergence of Markovian stochastic algorithms with rapidly decreasing ergodicity rates. *Stochastics and Stochastic Reports*, **65**(3), 177–228.