

<https://arxiv.org/pdf/1406.2661>

# GAN: Generative Adversarial Nets

2024.11.18 | 최민혜

# Contents

0 Abstract

1 Introduction

2 Adversarial Nets

3 Theoretical Results

4 Experiments

5 Advatanges and  
Disadvantages

6 Conclustions and future work

# 0. Abstract

## Abstract

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.

- 'Adversarial' 적대적인 : 두 개의 네트워크를 사용
- generative model (생성자),  $G$ 
  - : 분별모델을 상대로 완벽한 속임수를 수행
  - : training data의 분포를 모사
- discriminative model (판별자),  $D$ 
  - : 실제 데이터와 생성 모델이 만들어낸 데이터를 구별
  - : sample data가  $G$ 로부터 나온 데이터가 아닌 실제 training data로부터 나온 데이터일 확률을 추정

→ minimax two-player game

# 1. Introduction

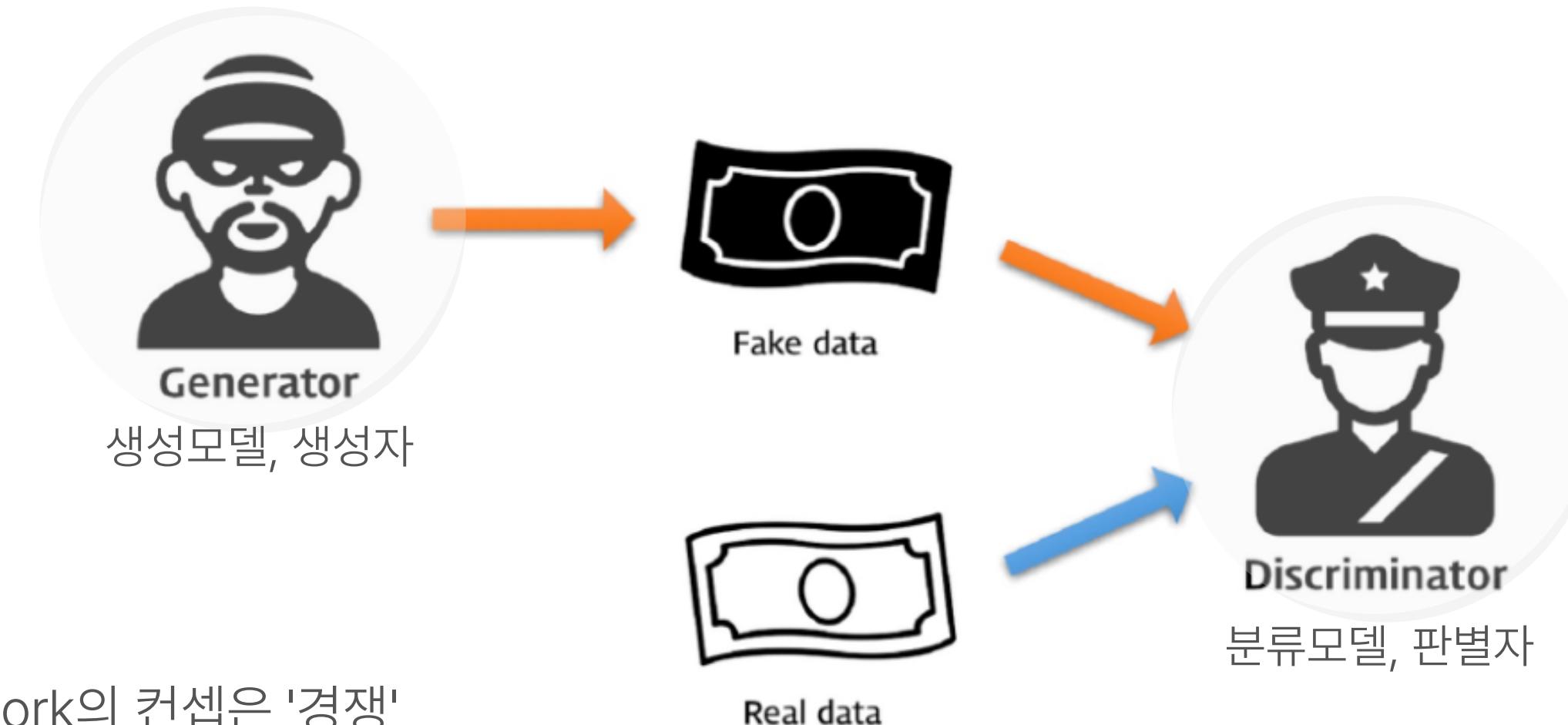
## 1 Introduction

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>

- 2015년 - deep learning 을 활용한 discriminative model이 많이 성장하던 시기
- Deep learning 기반의 generative model은 몇 가지 어려움이 존재
- 이러한 상황에서 어려움들을 적절히 회피하는 새로운 모델 GAN을 제안

# 1. Introduction

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.



- adversarial nets framework의 컨셉은 '경쟁'

# 2. Related Work

## 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

목적함수

$$\min_G \max_D V(D, G) = \underbrace{\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})]}_{\text{첫 번째 항}} + \underbrace{\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]}_{\text{두 번째 항}}.$$

- 첫 번째 항 : real data  $\mathbf{x}$ 를 discriminator에 넣었을 때 나오는 결과를 log 취했을 때 얻는 기댓값
- 두 번째 항 : fake data  $\mathbf{z}$ 를 generator에 넣었을 때 나오는 결과를 discriminator에 넣었을 때 그 결과를  $\log(1 - \text{결과})$  했을 때 얻는 기댓값

학습 초반

$$D(G(\mathbf{z})) = 0 \text{에 가까움}$$

$\mathbf{z}$ 로부터 생성해낸 fake 이미지가  
가짜라고 D가 판별할 수 있음



학습 진행

$$D(G(\mathbf{z})) = 1 \text{에 가까워짐}$$

$\mathbf{z}$ 로부터 생성해낸 fake 이미지가  
D는 진짜라고 판별해 버리게 됨

# 3. Adversarial nets

D 관점

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

G 관점

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

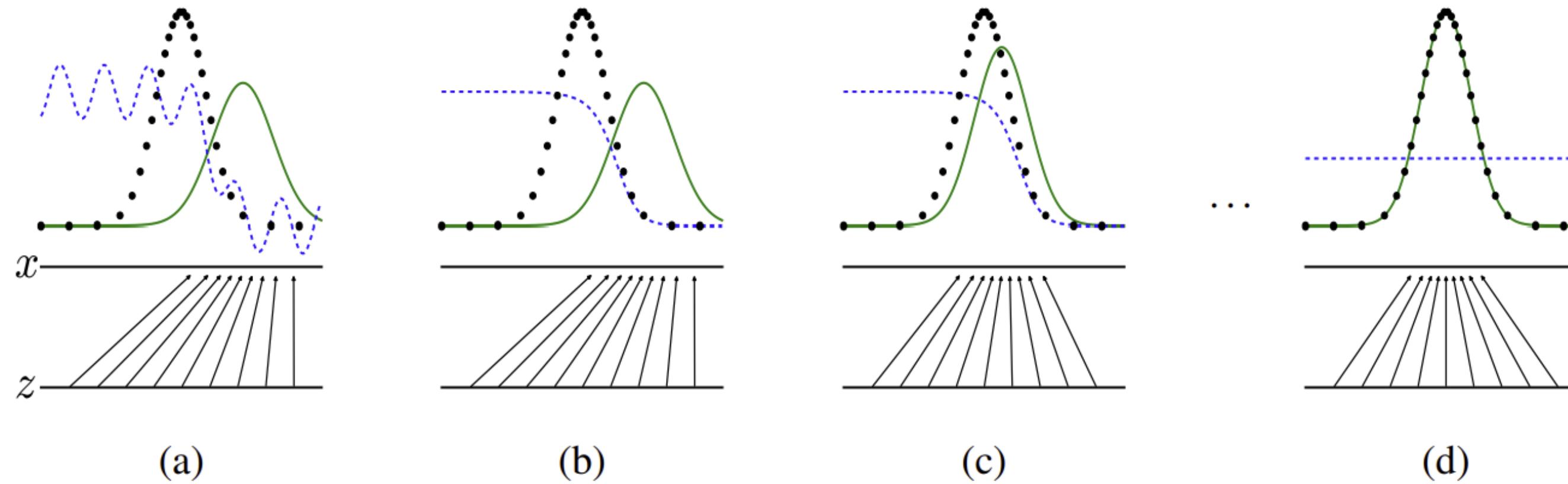
G의 성능에 의해 결정될 수 있는 항이 아님

→ D는  $V(D, G)$  를 최대화 , G는  $V(D, G)$  를 최소화하기 위해 학습

→ D와 G를  $V(D, G)$ 를 가지는 two-player minmax game이라 표현

# 3. Adversarial nets

- 파란색 점선: discriminative distribution
- 검은색 점선: data generating distribution(real)
- 녹색 실선: generative distribution(fake)



- (a): 학습 초기 - real과 fake의 분포가 전혀 다름
- (b): a처럼 불분명하게 확률을 판단하지 않고, 흔들리지 않고 real과 fake를 분명하게 판별해냄(D의 성능 향상)
- (c): 어느 정도 D가 학습이 이루어지면, G는 실제 데이터의 분포를 모사하며 D가 구별하기 힘든 방향으로 학습함
- (d): 이 과정의 반복의 결과로 real과 fake의 분포가 거의 비슷해져 구분할 수 없을 만큼 G가 학습을 함.  
결국 둘을 구분할 수 없게 되어 확률을 1/2로 계산하게 됨

# 4. Theoretical Results

---

**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**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
- Update the discriminator by ascending its stochastic gradient:

        K step

D

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

**end for**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Update the generator by descending its stochastic gradient:

        1 step

G

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

**end for**

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

# 4. Theoretical Results

## 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*

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \quad (2)$$

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

$$\begin{aligned} V(G, D) &= \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{z}} p_{\mathbf{z}}(\mathbf{z}) \log(1 - D(g(\mathbf{z}))) d\mathbf{z} \\ &= \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x} \end{aligned} \quad (3)$$

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.  $\square$

Note that the training objective for  $D$  can be interpreted as maximizing the log-likelihood for estimating the conditional probability  $P(Y = y|\mathbf{x})$ , where  $Y$  indicates whether  $\mathbf{x}$  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:

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

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

*Proof.* For  $p_g = p_{\text{data}}$ ,  $D_G^*(\mathbf{x}) = \frac{1}{2}$ , (consider Eq. 2). Hence, by inspecting Eq. 4 at  $D_G^*(\mathbf{x}) = \frac{1}{2}$ , 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  $p_g = p_{\text{data}}$ , observe that

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

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

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

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:

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

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  $p_g = p_{\text{data}}$ , i.e., the generative model perfectly replicating the data generating process.  $\square$

# 4. Theoretical Results

## 4.2 Convergence of Algorithm 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*

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

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

*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.  $\square$

In practice, adversarial nets represent a limited family of  $p_g$  distributions via the function  $G(\mathbf{z}; \theta_g)$ , 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

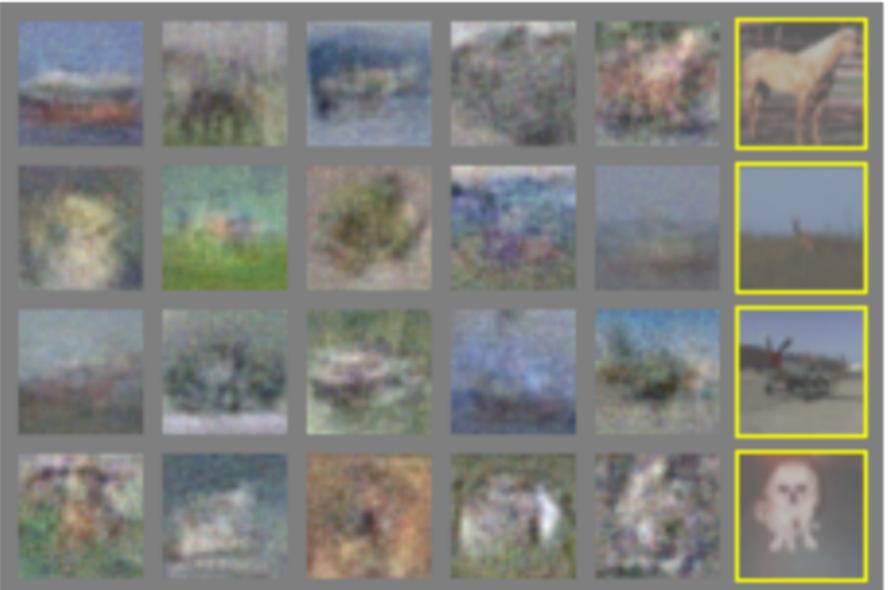
Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 66$
Stacked CAE [3]	$121 \pm 1.6$	$2110 \pm 50$
Deep GSN [6]	$214 \pm 1.1$	$1890 \pm 29$
Adversarial nets	<b><math>225 \pm 2</math></b>	<b><math>2057 \pm 26</math></b>



a)



b)



c)



d)

# 6. Advantaged and Disadvantages

## 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  $p_g(\mathbf{x})$ , 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  $\mathbf{z}$  to the same value of  $\mathbf{x}$  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.

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.

- 단점
  - $p_g(\mathbf{x})$  가 명시적으로 표현되지 않음
  - D와 G가 균형을 잘 맞춰 성능이 향상되어야 함
- 장점
  - Markov chains이 전혀 필요없음
  - 학습 중 어떠한 inference가 필요없음
  - 다양한 함수들이 모델에 접목될 수 있음

# 7. Conclusions and future work

## 7 Conclusions and future work

This framework admits many straightforward extensions:

1. A *conditional generative* model  $p(\mathbf{x} \mid \mathbf{c})$  can be obtained by adding  $\mathbf{c}$  as input to both  $G$  and  $D$ .
2. *Learned approximate inference* can be performed by training an auxiliary network to predict  $\mathbf{z}$  given  $\mathbf{x}$ . 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  $p(\mathbf{x}_S \mid \mathbf{x}_{\bar{S}})$  where  $S$  is a subset of the indices of  $\mathbf{x}$  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  $\mathbf{z}$  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.

- conditional generative model로 발전시킬 수 있음
- parameters를 공유하는 conditional model를 학습함으로써 다른 conditional models을 근사적으로 모델링할 수 있음
- Learned approximate inference는 주어진  $\mathbf{x}$ 를 예측하여 수행될 수 있음
- 효율성 개선:  $G, D$ 를 조정하는 더 나은 방법이나 학습하는 동안 sample  $\mathbf{z}$ 에 대한 더 나은 분포를 결정함으로써 학습의 속도를 높일 수 있음
- Semi-supervised learning: 제한된 레이블이 있는 데이터를 사용할 수 있을 때, classifiers의 성능을 향상시킬 수 있음

# 감사합니다