

# 0705 Paper Review

---

21 박연수

# Index

- 1. GAN
    - Adversarial nets
    - Theoretical Results
    - Experiments
    - Advantages and disadvantages
    - Conclusions and future work
-

## Statistics &gt; Machine Learning

[Submitted on 10 Jun 2014]

# Generative Adversarial Networks

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

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  $1/2$  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.

Subjects: **Machine Learning (stat.ML)**; Machine Learning (cs.LG)

Cite as: arXiv:1406.2661 [stat.ML]  
(or arXiv:1406.2661v1 [stat.ML] for this version)  
<https://doi.org/10.48550/arXiv.1406.2661>

## Submission history

From: Ian Goodfellow [view email]  
[v1] Tue, 10 Jun 2014 18:58:17 UTC (1,257 KB)

## Bibliographic Tools

Code, Data, Media

Demos

Related Papers

About arXivLabs

## Bibliographic and Citation Tools

☐ Bibliographic Explorer (What is the Explorer?)

☐ Litmaps (What is Litmaps?)

☐ scite Smart Citations (What are Smart Citations?)

Which authors of this paper are endorsers? | Disable MathJax (What is MathJax?)

## Access Paper:

- [View PDF](#)
- [TeX Source](#)
- [Other Formats](#)

[view license](#)

Current browse context:

stat.ML

&lt; prev | next &gt;

new | recent | 2014-06

Change to browse by:

cs

cs.LG

stat

## References & Citations

- [NASA ADS](#)
- [Google Scholar](#)
- [Semantic Scholar](#)

59 blog links (what is this?)

Export BibTeX Citation

## Bookmark



Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.  
Generative Adversarial Networks.  
In NIPS, 2014  
[1406.2661 \(arxiv.org\)](https://arxiv.org/abs/1406.2661)

## 1. GAN

### Adversarial nets

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

- G is generator mapping noise z to data space.
  - D is discriminator represents the probability that x came from the data rather than pg
  - G to minimize  $\log(1 - D(G(z)))$
  - D to maximize  $\log(D(x))$
  - D and G play the following two-player minimax game with value function  $V(G, D)$
-

## 1. GAN

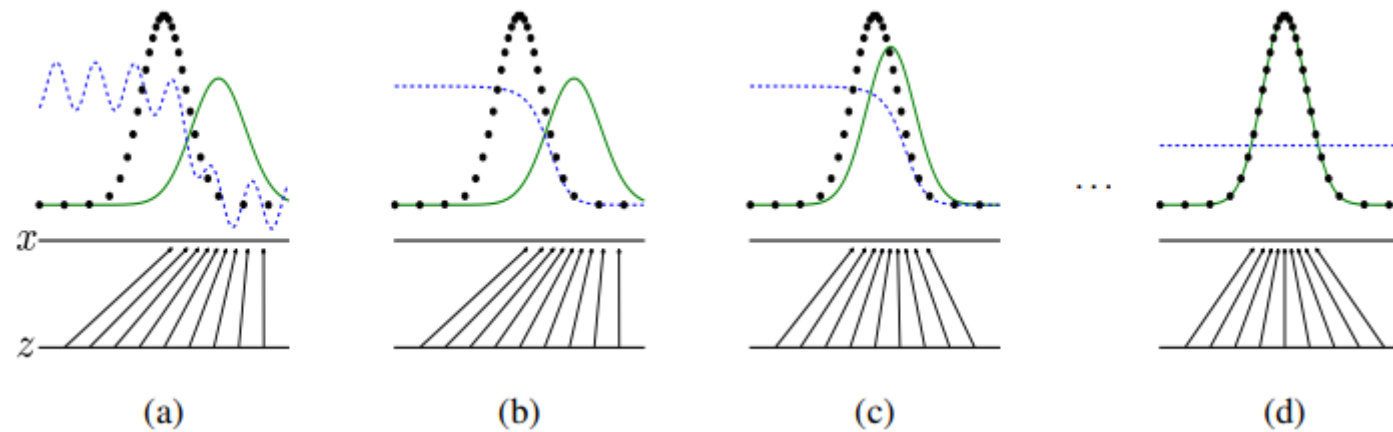
### Adversarial nets

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

- Optimizing D to completion in the inner loop of training is computationally prohibitive
  - k steps of optimizing D and one step of optimizing G. -> D being maintained near its optimal solution, so long as G changes slowly enough.
  - G to minimize  $\log(1 - D(G(\mathbf{z})))$  VS G to maximize  $\log D(G(\mathbf{z}))$ 
    - In practice, equation may not provide sufficient gradient for G to learn well.
    - when G is poor, D can reject samples with high confidence In this case,  $\log(1 - D(G(\mathbf{z})))$  saturates.
    - We can train G to maximize  $\log D(G(\mathbf{z}))$
-

## 1. GAN

### Adversarial nets



- $D$  : blue, dashed line
- $G/p_g$  : green, solid line
- Data : black, dotted line
- Arrows :  $x = G(z)$

- (a) Consider an adversarial pair near convergence:  $p_g$  is similar to  $p_{data}$  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
$$D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$
- (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  $p_g = p_{data}$ . The discriminator is unable to differentiate between the two distributions, i.e.  $D(x) = 1/2$

## 1. GAN

### 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  $\{z^{(1)}, \dots, z^{(m)}\}$  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:

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

**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:

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

**end for**

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

---

## 1. GAN

### Theoretical Results

- Why KL divergence appears in this paper?

$$D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

$$C(G) = \max_D V(G, D)$$

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

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

$$= \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right]$$

$$KL(B \parallel A) = \mathbb{E}_{x \sim B} \left[ \log \frac{B(x)}{A(x)} \right]$$

let A as  $p_g$  and let B as  $p_{data}$

$$C(G) = -\log 4 + \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{2p_{data}(x)}{p_{data}(x) + p_g} \right] + \mathbb{E}_{x \sim p_z} \left[ \log \frac{2p_g(x)}{p_{data}(x) + p_g} \right]$$

$$= -\log 4 + KL(p_{data} \parallel \frac{p_{data} + p_g}{2}) + KL(p_g \parallel \frac{p_{data} + p_g}{2})$$

$$JSD(p \parallel q) = \frac{1}{2} (KL(p \parallel M) + KL(q \parallel M))$$

$$\text{where, } M = \frac{1}{2}(p + q)$$

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



## 1. GAN

### Theoretical Results

- Why KL divergence appears in this paper?

$$D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

$$C(G) = \max_D V(G, D)$$

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

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

$$= \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right]$$

$$KL(B \parallel A) = \mathbb{E}_{x \sim B} \left[ \log \frac{B(x)}{A(x)} \right]$$

let A as  $p_g$  and let B as  $p_{data}$

$$C(G) = -\log 4 + \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{2p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g} \right] + \mathbb{E}_{\mathbf{x} \sim p_z} \left[ \log \frac{2p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g} \right]$$

$$= -\log 4 + KL(p_{data} \parallel \frac{p_{data} + p_g}{2}) + KL(p_g \parallel \frac{p_{data} + p_g}{2})$$

$$JSD(p \parallel q) = \frac{1}{2} (KL(p \parallel M) + JSD(q \parallel M))$$

$$\text{where, } M = \frac{1}{2}(p + q)$$

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

## 1. GAN

### Theoretical Results

- Why KL divergence appears in this paper?

$$D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

$$C(G) = \max_D V(G, D)$$

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

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

$$= \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right]$$

$$KL(B \parallel A) = \mathbb{E}_{x \sim B} \left[ \log \frac{B(x)}{A(x)} \right]$$

let A as  $p_g$  and let B as  $p_{data}$

$$C(G) = -\log 4 + \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{2p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g} \right] + \mathbb{E}_{\mathbf{x} \sim p_z} \left[ \log \frac{2p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g} \right]$$

$$= -\log 4 + KL(p_{data} \parallel \frac{p_{data} + p_g}{2}) + KL(p_g \parallel \frac{p_{data} + p_g}{2})$$

$$JSD(p \parallel q) = \frac{1}{2} (KL(p \parallel M) + JSD(q \parallel M))$$

$$\text{where, } M = \frac{1}{2}(p + q)$$

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

## 1. GAN

### Advantages and disadvantages

Advantages	Disadvantages
<ol style="list-style-type: none"><li>1. Markov chains are never needed</li><li>2. no inference is needed</li><li>3. wide variety of functions can be incorporated into the model</li></ol>	<ol style="list-style-type: none"><li>1. there is no explicit representation of <math>p_g(x)</math></li><li>2. D must be synchronized well with G during training</li></ol>

## 1. GAN

### Conclusions and future work

- A conditional generative model
  - Learned approximate inference
  - Semi-supervised learning
  - Etc ...
- 
- This paper has demonstrated the viability of the adversarial modeling framework, suggesting that these research directions could prove useful.
-

Thank you

---