

# Generative Adversarial Nets

He Zean <sup>1</sup>

<sup>1</sup>Dept. of Computer Science and Engineering, SUSTech

Nov. 9, 2022



# Outline

- 1 Introduction
- 2 Method
- 3 Experiments
- 4 Analysis
- 5 Code Implementation

# Outline

## 1 Introduction

- Background
- Motivation

## 2 Method

## 3 Experiments

## 4 Analysis

## 5 Code Implementation



# Introduction: Overview

# Introduction: Overview

# Introduction: Overview

- Generative Adversarial Nets (GAN) was first introduced by Goodfellow *et al.* in 2014 [1].

# Introduction: Overview

- Generative Adversarial Nets (GAN) was first introduced by Goodfellow *et al.* in 2014 [1].
- Popular Applications:

# Introduction: Overview

- Generative Adversarial Nets (GAN) was first introduced by Goodfellow *et al.* in 2014 [1].
- Popular Applications:
  - Image Editing



# Introduction: Overview

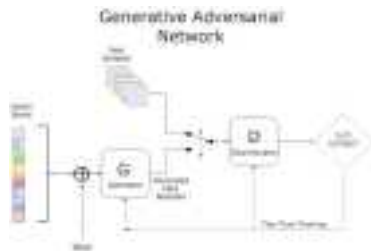
- Generative Adversarial Nets (GAN) was first introduced by Goodfellow *et al.* in 2014 [1].
- Popular Applications:
  - Image Editing
  - Data Generating

# Introduction: Overview

- Generative Adversarial Nets (GAN) was first introduced by Goodfellow *et al.* in 2014 [1].
- Popular Applications:
  - Image Editing
  - Data Generating
  - Cross-Domain Image Translation [2]

# Introduction: Overview

- Generative Adversarial Nets (GAN) was first introduced by Goodfellow *et al.* in 2014 [1].
- Popular Applications:
  - Image Editing
  - Data Generating
  - Cross-Domain Image Translation [2]



Retrieved from  
<https://paperswithcode.com/method/gan>.

# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)

# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)
  - The partition function & gradient are hard to calculate.

# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)
  - The partition function & gradient are hard to calculate.
  - Mixing is a problem for learning algorithms that rely on MCMC.

# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)
  - The partition function & gradient are hard to calculate.
  - Mixing is a problem for learning algorithms that rely on MCMC.
- Deep directed graphical models (ex. DBNs, hybrid)

# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)
  - The partition function & gradient are hard to calculate.
  - Mixing is a problem for learning algorithms that rely on MCMC.
- Deep directed graphical models (ex. DBNs, hybrid)
  - Computational difficulties: a fast approximate layer-wise training criterion exists.



# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)
  - The partition function & gradient are hard to calculate.
  - Mixing is a problem for learning algorithms that rely on MCMC.
- Deep directed graphical models (ex. DBNs, hybrid)
  - Computational difficulties: a fast approximate layer-wise training criterion exists.
- Generative autoencoders (ex. NCE, GSNs)

# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)
  - The partition function & gradient are hard to calculate.
  - Mixing is a problem for learning algorithms that rely on MCMC.
- Deep directed graphical models (ex. DBNs, hybrid)
  - Computational difficulties: a fast approximate layer-wise training criterion exists.
- Generative autoencoders (ex. NCE, GSNs)
  - Fixed noise distribution, learning slows dramatically if the distribution of even a small subset of data has been learnt.

# Introduction: Background

- Deep undirected graphical models (ex. RBMs, DBMs)
  - The partition function & gradient are hard to calculate.
  - Mixing is a problem for learning algorithms that rely on MCMC.
- Deep directed graphical models (ex. DBNs, hybrid)
  - Computational difficulties: a fast approximate layer-wise training criterion exists.
- Generative autoencoders (ex. NCE, GSNs)
  - Fixed noise distribution, learning slows dramatically if the distribution of even a small subset of data has been learnt.
  - Performance problem: requires feedback loops, less able to leverage piecewise linear units.

# Introduction: Motivation

- We want to construct a generative model that have lower computational costs, comparing with existing generative models.

# Introduction: Motivation

- We want to construct a generative model that have lower computational costs, comparing with existing generative models.
- We want to sidestep difficulties such as intractable probabilistic computations and leveraging piecewise linear units.

# Outline

- 1 Introduction
- 2 Method**
- 3 Experiments
- 4 Analysis
- 5 Code Implementation

# Method

# Method

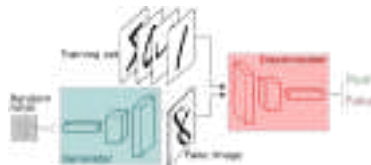


- **Generator:** input a random noise, the generator "convert" it to a fake image, and try to cheat the discriminator.

- **Generator:** input a random noise, the generator "convert" it to a fake image, and try to cheat the discriminator.
- **Discriminator:** input an image (n-dim vector), output a label: the image comes from training data (real) or generator (fake). To be practical in training, we directly use the estimated possibility.

# Method

- **Generator:** input a random noise, the generator "convert" it to a fake image, and try to cheat the discriminator.
- **Discriminator:** input an image (n-dim vector), output a label: the image comes from training data (real) or generator (fake). To be practical in training, we directly use the estimated possibility.



Retrieved from  
<https://zhuanlan.zhihu.com/p/33752313>.

Example: A multilayer perceptron - multilayer perceptron GAN

w.r.t.:

Example: A multilayer perceptron - multilayer perceptron GAN

$$G(\mathbf{z}; \theta_g), \mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$$

w.r.t.:

Example: A multilayer perceptron - multilayer perceptron GAN

$$G(\mathbf{z}; \theta_g), \mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$$

$$D(\mathbf{x}; \theta_d) = p(\mathbf{x} \in \text{data})$$

w.r.t.:

Example: A multilayer perceptron - multilayer perceptron GAN

$$G(\mathbf{z}; \theta_g), \mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$$

$$D(\mathbf{x}; \theta_d) = p(\mathbf{x} \in \text{data})$$

w.r.t.:

- $G(\mathbf{z}; \theta_g)$  is the function representation of generator (multilayer perceptron) with parameters  $\theta_g$ , differentiable.  $D(\mathbf{x}; \theta_d)$  similars.

Example: A multilayer perceptron - multilayer perceptron GAN

$$G(\mathbf{z}; \theta_g), \mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$$

$$D(\mathbf{x}; \theta_d) = p(\mathbf{x} \in \text{data})$$

w.r.t.:

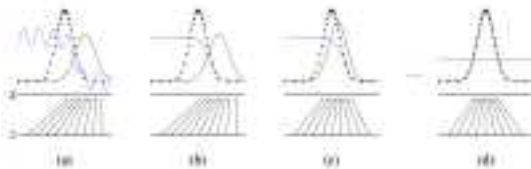
- $G(\mathbf{z}; \theta_g)$  is the function representation of generator (multilayer perceptron) with parameters  $\theta_g$ , differentiable.  $D(\mathbf{x}; \theta_d)$  similars.
- $p_{\mathbf{z}}(\mathbf{z})$  is the prior of input noise



# Method: Training Process

POV: Minimax gaming

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



# Method: Training Process

- Global Optimality of  $p_g = p_{data}$

# Method: Training Process

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

# Method: Training Process

- $d\_loss = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log (1 - D(G(\mathbf{z})))]$ , cross-entropy.

---

<sup>1</sup>Early in learning, when G is poor, rather than training G to minimize  $[\log (1 - D(G(\mathbf{z})))]$ , we can train G to maximize  $[\log (D(G(\mathbf{z})))]$ .

# Method: Training Process

- $d\_loss = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log (1 - D(G(\mathbf{z})))]$ , cross-entropy.
- $g\_loss = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log (1 - D(G(\mathbf{z})))]$ <sup>1</sup>

---

<sup>1</sup>Early in learning, when G is poor, rather than training G to minimize  $[\log (1 - D(G(\mathbf{z})))]$ , we can train G to maximize  $[\log (D(G(\mathbf{z})))]$ .

# Method: Training Process

**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  $p_z(\mathbf{z})$ .
- 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:

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

**end for**

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

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

**end for**

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

# Outline

- 1 Introduction
- 2 Method
- 3 Experiments**
- 4 Analysis
- 5 Code Implementation

# Experiments

- Datasets

Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 68$
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>



# Experiments

- Datasets
  - MNIST, TFD, CIFAR-10

Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 68$
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>

# Experiments

- Datasets
  - MNIST, TFD, CIFAR-10
- NNs' design will be mentioned in next section

Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 68$
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>

# Experiments

- Datasets
  - MNIST, TFD, CIFAR-10
- NNs' design will be mentioned in next section
- $p_g$  estimation

Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 68$
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>

# Experiments

- Datasets
  - MNIST, TFD, CIFAR-10
- NNs' design will be mentioned in next section
- $p_g$  estimation
  - Gaussian Parzen window

Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 68$
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>

# Experiments

- Datasets
  - MNIST, TFD, CIFAR-10
- NNs' design will be mentioned in next section
- $p_g$  estimation
  - Gaussian Parzen window
  - Cross validation  $\rightarrow \sigma$

Model	MNIST	TFD
DBN [3]	$138 \pm 2$	$1909 \pm 68$
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>

# Outline

- 1 Introduction
- 2 Method
- 3 Experiments
- 4 Analysis**
- 5 Code Implementation

- Advantages

- Advantages
  - Less computational cost



- Advantages
  - Less computational cost
  - Model can accept multiple types of functions (differentiable)

- Advantages
  - Less computational cost
  - Model can accept multiple types of functions (differentiable)
  - Is able to represent very sharp, even degenerate distributions

- Advantages
  - Less computational cost
  - Model can accept multiple types of functions (differentiable)
  - Is able to represent very sharp, even degenerate distributions
- Disadvantages

- Advantages

- Less computational cost
- Model can accept multiple types of functions (differentiable)
- Is able to represent very sharp, even degenerate distributions

- Disadvantages

- Required synchronization between G and D (to avoid mode collapse)

- Advantages

- Less computational cost
- Model can accept multiple types of functions (differentiable)
- Is able to represent very sharp, even degenerate distributions

- Disadvantages

- Required synchronization between G and D (to avoid mode collapse)
- Highly sensitive to the hyperparameter selections

# Outline

- 1 Introduction
- 2 Method
- 3 Experiments
- 4 Analysis
- 5 Code Implementation**
  - Data Format
  - Descriptions
  - Result

# Input Format

```
data_loader = torch.utils.data.DataLoader(  
    datasets.MNIST(  
        root=f'../data/mnist',  
        train=True,  
        download=True,  
        transform=transforms.Compose(  
            [transforms.Resize(img_size), transforms.ToTensor(), transforms.Normalize([0.5], [0.5])] )  
    ),  
    batch_size=batch_size,  
    shuffle=True,  
)
```

- MNIST dataset provides 60,000 samples, they are split into  $938 \times 64$ -sized mini-batch for each epoch

# Input Format

```
data_loader = torch.utils.data.DataLoader(  
    datasets.MNIST(  
        root=f'../data/mnist',  
        train=True,  
        download=True,  
        transform=transforms.Compose(  
            [transforms.Resize(img_size), transforms.ToTensor(), transforms.Normalize([0.5], [0.5])] )  
    ),  
    batch_size=batch_size,  
    shuffle=True,  
)
```

- MNIST dataset provides 60,000 samples, they are split into  $938 \times 64$ -sized mini-batch for each epoch
- Each sample is a  $28 \times 28$  pixel, 1 channel (grayscale) image



# Input Format

```
data_loader = torch.utils.data.DataLoader(
    datasets.MNIST(
        root='./data/mnist',
        train=True,
        download=True,
        transform=transforms.Compose(
            [transforms.Resize(img_size), transforms.ToTensor(), transforms.Normalize([0.5], [0.5])]
        ),
    batch_size=batch_size,
    shuffle=True,
```

- MNIST dataset provides 60,000 samples, they are split into  $938 \times 64$ -sized mini-batch for each epoch
- Each sample is a  $28 \times 28$  pixel, 1 channel (grayscale) image
- ToTensor() + Normalize([0.5], [0.5]): each pixel's value is mapped to the range of  $[-1, 1]$

# Input Format (Cont.)

```
# sample noise as generator input
z = Variable(Tensor(np.random.normal(0, 1, [args.shape[0], opt.latent_dim])))
```

- For each mini-batch, it generates 64 (batch size) 100-dim (latent dimension) Gaussian noises ( $\mu = 0, \sigma = 1$ ) as noise prior  $p_g(\mathbf{z})$

- Optimizer: Adam

- Optimizer: Adam
- Adversarial loss: Binary Classification Error Loss

- Optimizer: Adam
- Adversarial loss: Binary Classification Error Loss

```
# Adversarial ground truths  
neg_label = keras.utils.to_categorical(np.zeros((1, 10)), 10, dtype='float32')  
neg_label = keras.utils.to_categorical(np.zeros((1, 10)), 10, dtype='float32')
```

# Output Format

```
print(
    "\tEpoch %d/%d | Batch %d/%d | D loss: %f | G loss: %f"
    % (epoch, opt.n_epochs, i, len(dataloader), d_loss.item(), g_loss.item())
)

batches_done = epoch * len(dataloader) + 1
if batches_done % opt.sample_interval == 0:
    save_image(gen_imgs.data[0], "%s/%d.png" % (batches_done, row*5, serialize=True))
```

- For each mini-batch, it uses BCELoss to calculate  $d\_loss$  (loss for discriminator) and  $g\_loss$  (loss for generator)

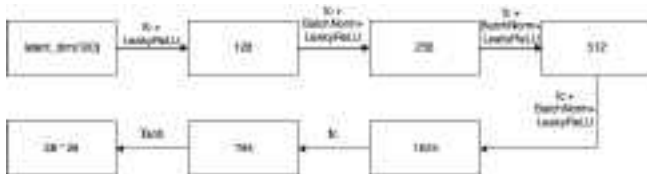
# Output Format

```
print(
    "[Epoch %d/%d] [Batch %d/%d] [D loss: %f] [G loss: %f]"
    % (epoch, opt.n_epochs, 1, len(dataloader), d_loss.item(), g_loss.item())

batches_done = epoch + len(dataloader) + 1
if batches_done % opt.sample_interval == 0:
    save_image(gan_imgs.data[:25], "images/%d.png" % batches_done, row=0, normalize=True)
```

- For each mini-batch, it uses BCELoss to calculate  $d\_loss$  (loss for discriminator) and  $g\_loss$  (loss for generator)
- For every 400 batches (sample interval), it saves the first 25 generated images ( $28 \times 28$  pixel, 1 channel grayscale)

# Descriptions



Generator NN Structure



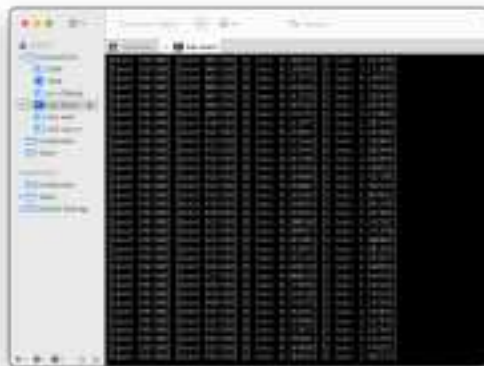
# Descriptions



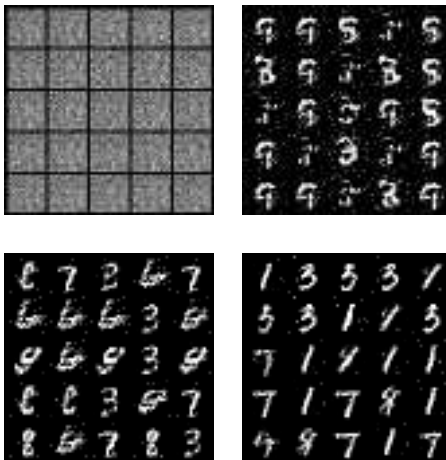
Discriminator NN Structure

# Running Result

Ran on server with GeForce RTX 2080 Ti  
200 epochs, batch size = 64, learning rate = 0.0002



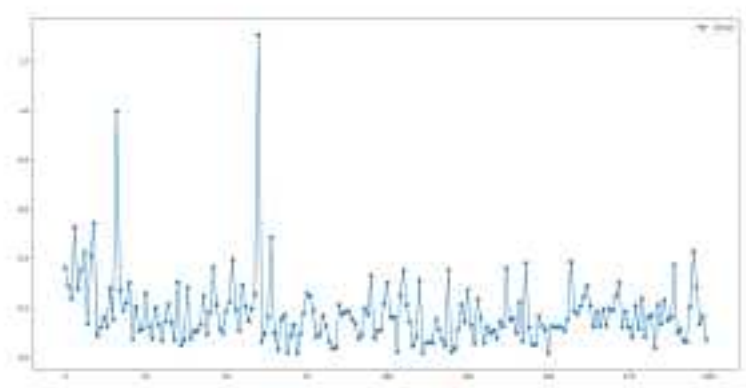
# Running Result



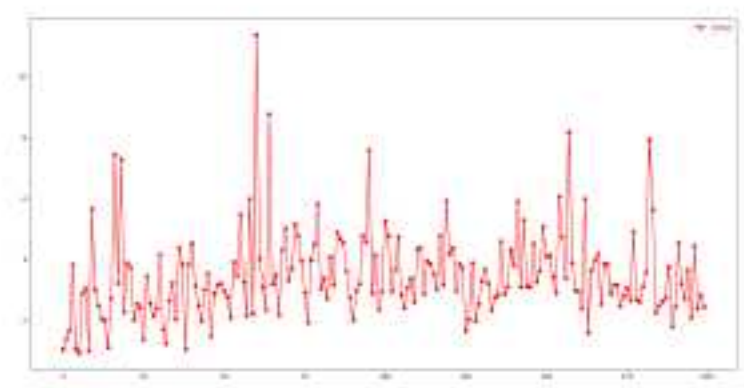
Generated Images for MNIST (patch #0, #46800, #93600, and # 180000)



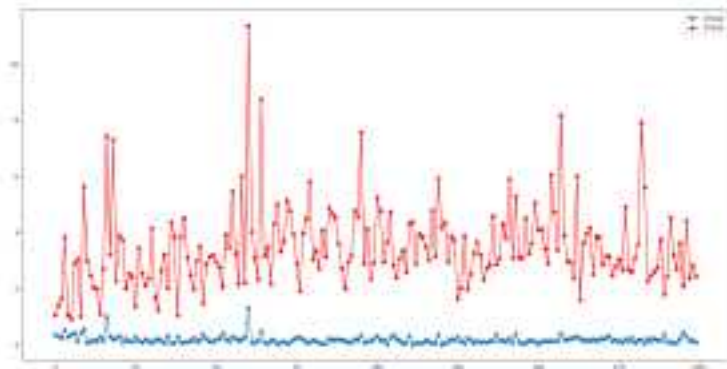
# Running Result



# Running Result



# Running Result



# References I

- [1] Ian J. Goodfellow **and others**. “Generative Adversarial Nets”. *in Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada: byeditor Zoubin Ghahramani and others*. 2014, **pages** 2672–2680.
- [2] Ming-Yu Liu, Thomas M. Breuel **and** Jan Kautz. “Unsupervised Image-to-Image Translation Networks”. *in Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA: byeditor Isabelle Guyon and others*. 2017, **pages** 700–708.