



# Generative Adversarial Nets

NIPS 2014

Ian J. Goodfellow et. al

2017.7.15

# 01 Adversarial nets

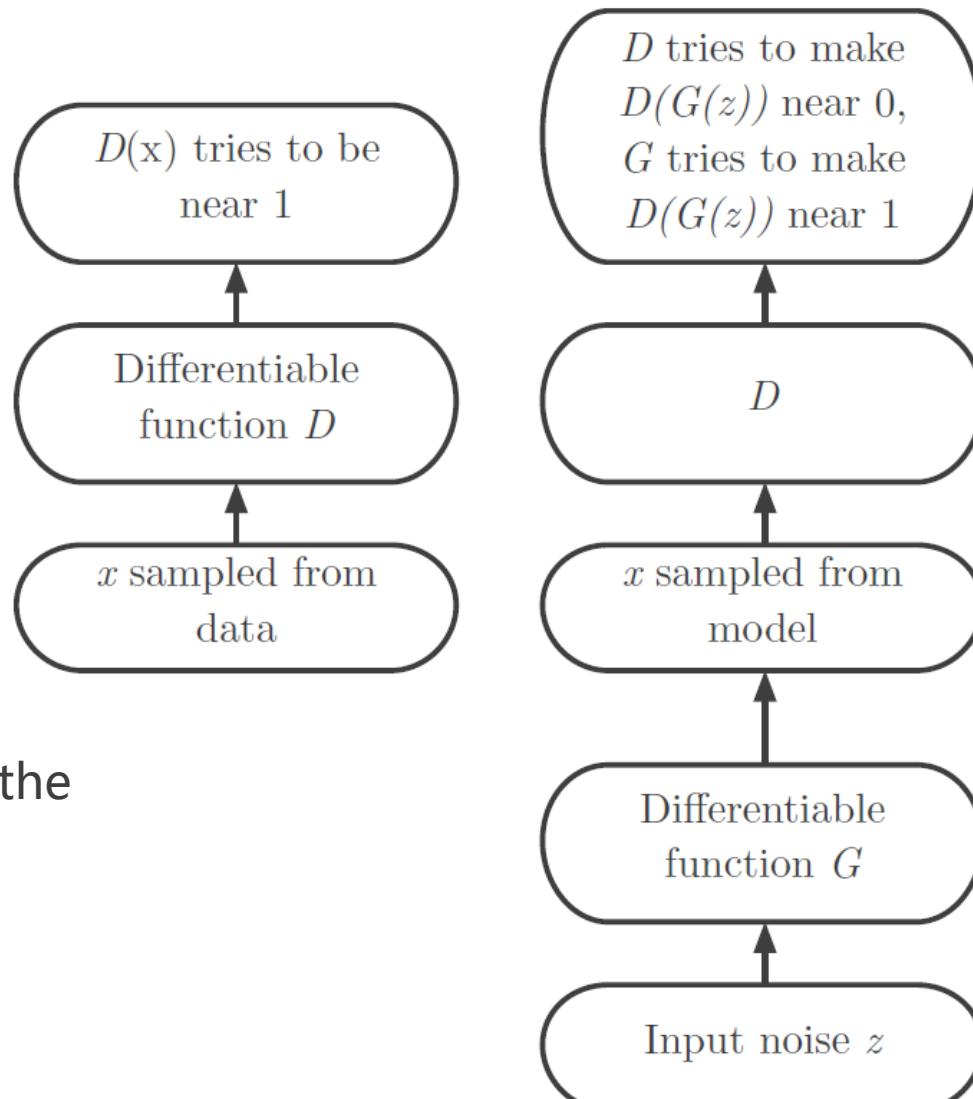


Input noise variables  $p_z(z)$

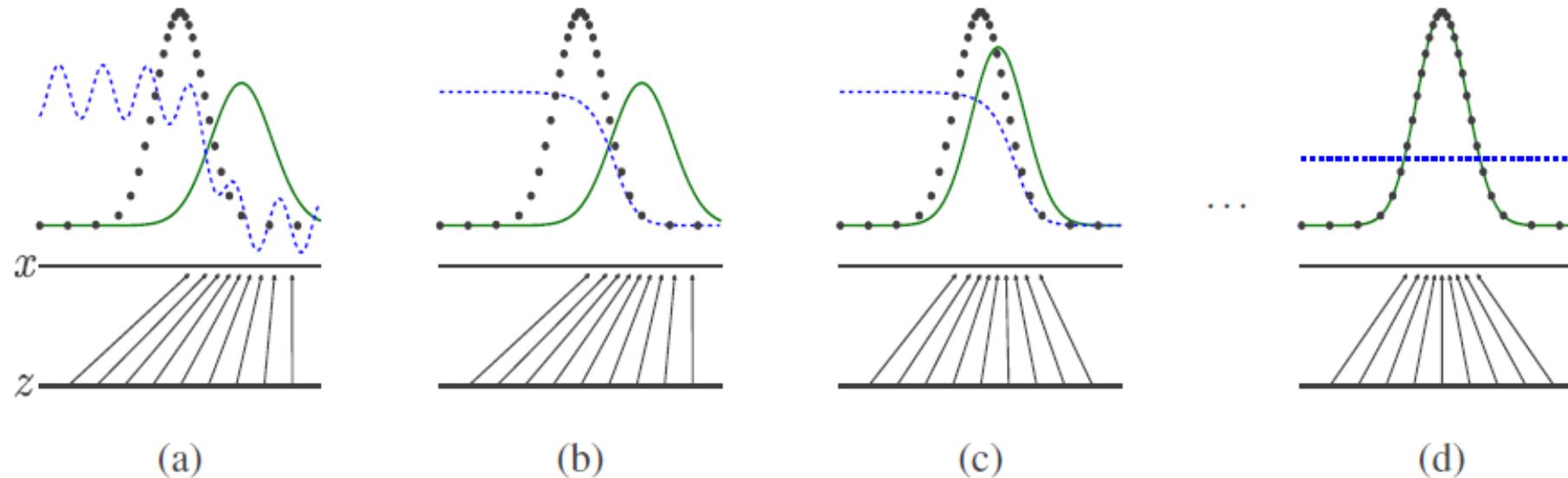
$G(z, \theta_g)$  represent a mapping to data space a  
 $G$  is a differentiable function represented by  
a multilayer perceptron with parameters  $\theta_g$ .

Multilayer perceptron  $D(x, \theta_d)$  outputs a single scalar.

$D(x)$  represents the probability that  $x$  came from the data rather than  $p_g$ .



# 01 Adversarial nets



D and G play the following two-player minimax game with value function  $V(G; D)$ :

# 02 Math



$$V(G, D) = E_{x \sim P_{data}}[\log D(x)] + E_{x \sim P_G}[\log(1 - D(x))] \quad G^* = \operatorname{argmin}_G (\max_D V(G, D))$$

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

$$V = E_{x \sim P_{data}}[\log D(x)] + E_{x \sim P_G}[\log(1 - D(x))]$$

$$= \int P_{data}(x) \log D(x) dx + \int p_G \log(1 - D(x)) dx$$

$$= \int [P_{data} \log D(x) + P_G(x) \log(1 - D(x))] dx$$

$$D^* \longrightarrow P_{data} \log D(x) + P_G(x) \log(1 - D(x))$$

## 02 Math



$$f(D) = a \log(D) + b \log(1 - D) \quad \frac{df(D)}{dD} = a \times \frac{1}{D} + b \log \times \frac{1}{1 - D} \times (-1) = 0$$

$$a \times \frac{1}{D^*} = b \times \frac{1}{1 - D^*} \quad D^* = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)}$$

$$\max V(G, D) = V(G, D^*)$$

$$= E_{x \sim P_{data}} \left[ \log \frac{P_{data}(x)}{P_{data}(x) + P_G(x)} \right] + E_{x \sim P_G} \left[ \log \frac{P_G(x)}{P_{data}(x) + P_G(x)} \right]$$

$$= \int P_{data}(x) \log \frac{\frac{1}{2} P_{data}}{P_{data}(x) + P_G(x)} dx + \int P_G(x) \log \frac{\frac{1}{2} P_G(x)}{P_{data} + P_G(x)} dx$$

# 02 Math



$$\begin{aligned} &= -2\log 2 + KL(P_{data} \mid\mid \frac{P_{data}(x) + P_G(x)}{2}) + KL(P_G(x) \mid\mid \frac{P_{data} + P_G(x)}{2}) \\ &= -2\log 2 + 2 \cdot JSD(p_{data} \mid\mid p_g) \end{aligned}$$

# 03 Framework



**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:

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

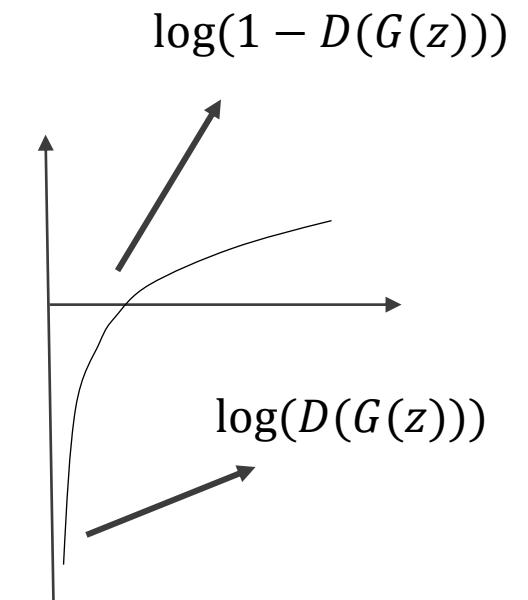
**end for**

- 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.





## 04 Cross-entropy cost function

$$C = \frac{1}{2n} \sum_x ||y(x) - a(x)||^2$$

$$C = \frac{1}{2}(y - a)^2$$

$$\frac{\partial C}{\partial w} = (a - y)\sigma(z)'x$$

$$\frac{\partial C}{\partial b} = (a - y)\sigma(z)'$$

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)]$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial C}{\partial w_j} = -\frac{1}{n} \sum_x \left[ \frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma(z)}{\partial w_j}$$
$$= -\frac{1}{n} \sum_x \left[ \frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \sigma(z)' x_j$$
$$\sigma(z)' = \sigma(z)(1 - \sigma(z))$$

$$= -\frac{1}{n} \sum_x \left[ \frac{\sigma(z)' x_j}{\sigma(z)(1 - \sigma(z))} (\sigma(z) - y) \right] = \frac{1}{n} \sum_x x_j (\sigma(z) - y)$$

$$\frac{\partial C}{\partial b} = \frac{1}{n} \sum_x (\sigma(z) - y)$$



## 04 Cross-entropy cost function

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial b} = \frac{\partial C}{\partial a} \sigma(z)' \frac{\partial (wx + b)}{\partial b} = \frac{\partial C}{\partial a} \sigma(z)' = \frac{\partial C}{\partial a} a(1 - a)$$

$$\frac{\partial C}{\partial b} = (a - y)\sigma(z)' \longrightarrow \frac{\partial C}{\partial b} = (a - y)$$

$$\frac{\partial C}{\partial a} a(1 - a) = a - y$$

$$C = -[y \ln a + (1 - y) \ln(1 - a)] + constant$$

# 05 Extra



Training processing

$$\theta^D \longrightarrow J^D \quad \theta^G \longrightarrow J^G$$

The discriminator's cost

$$J^D(\theta^D, \theta^G) = -\frac{1}{2} E_{x \sim p_{data}} \log(D(x)) - \frac{1}{2} E_z \log(1 - D(G(z)))$$

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

$$J^G = -J^D$$

# 05 Extra



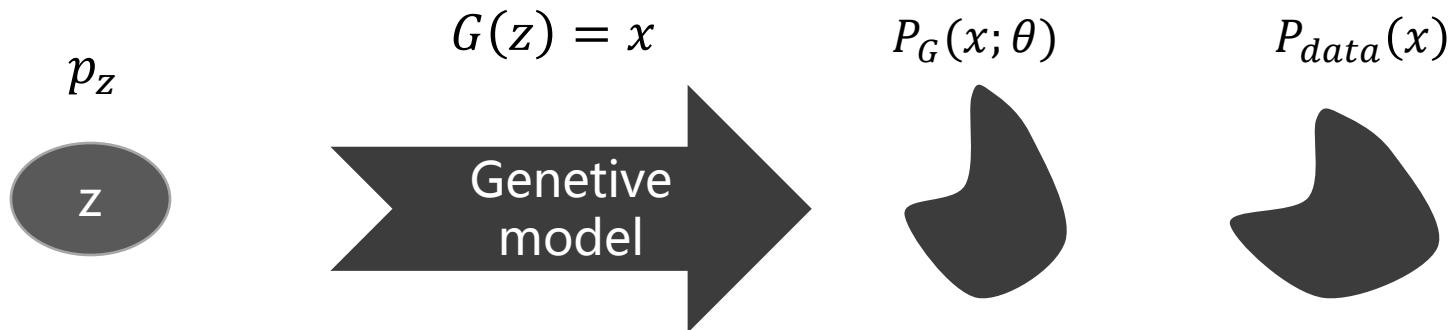
Because  $J^G$  is tied directly to  $J^D$ , we can summarize the entire game with a value function specifying the discriminator's payoff

$$V(\theta^D, \theta^G) = -J^D(\theta^D, \theta^G)$$

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

$$\theta^{G*} = \operatorname{argmin}_{\theta^G} \max_{\theta^D} V(\theta^D, \theta^G)$$

# 06 MLE



$$x^1, x^2, \dots, x^m$$

$$L = \prod_{i=1}^m P_G(x^i; \theta)$$

$$D_{KL}(P||Q) = \int_{-\infty}^{+\infty} p(x) \log \frac{p(x)}{q(x)}$$

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

$$\theta^* = \operatorname{argmax}_{\theta} \prod_{i=1}^m p_G(x^i, \theta) \Leftrightarrow \operatorname{argmax}_{\theta} \log \prod_{i=1}^m p_G(x^i, \theta)$$

# 06 MLE



$$= \operatorname{argmax}_{\theta} \sum_i^m \log P_G(x^i, \theta) \approx \operatorname{argmax}_{\theta} E_{x \sim P_{data}} [\log P_G(x; \theta)]$$

$$\Leftrightarrow \operatorname{argmax}_{\theta} \int P_{data}(x) \log P_G(x; \theta) dx - \int P_{data}(x) \log P_{data}(x) dx$$

$$= \operatorname{argmax}_{\theta} \int P_{data}(x) \log \frac{P_G(x; \theta)}{P_{data}(x)} dx$$

$$\operatorname{argmin}_{\theta} KL(P_{data} || P_G(x; \theta))$$

# 07 Advantages and disadvantages



There is no explicit representation of  $p_g$

D must be synchronized well with G during training

Markov chains are never needed, only backprop is used to obtain gradients

No inference is needed during learning

A wide variety of functions can be incorporated into the model