

Generative Adversarial Networks: Theory

Sargur N. Srihari

srihari@cedar.buffalo.edu

Topics in GAN

1. GAN Motivation

2. GAN Theory

1. Generator and Discriminator
2. Adversarial Modeling framework
3. GAN Objective and SGD
4. Theoretical Results
5. GAN Difficulties

→ Discussed Here

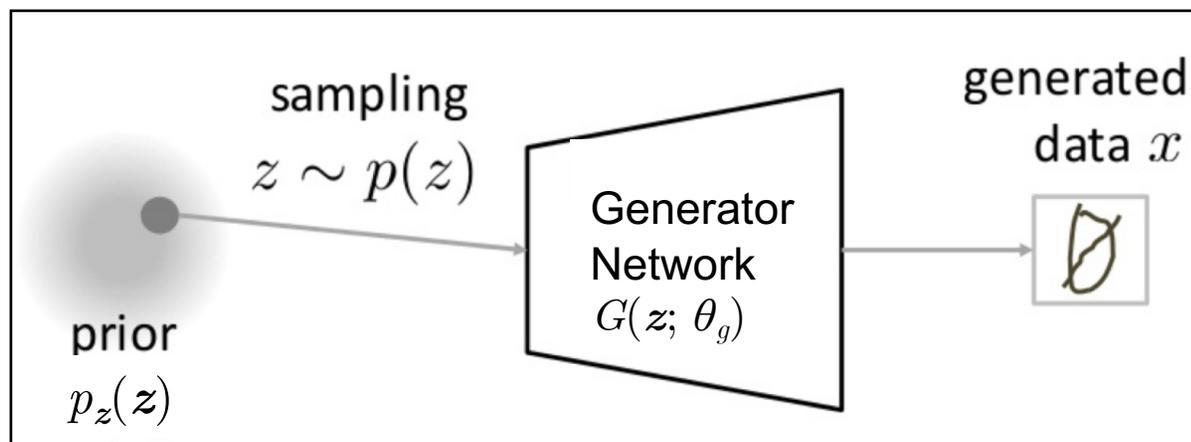
3. Mode Collapse

4. Wasserstein GAN

5. GAN Variants

Generator Network

- To learn the generator's distribution p_g over data \mathbf{x} , we define a prior on input noise variables $p_z(\mathbf{z})$
- Represent a mapping to data space as $G(\mathbf{z}; \theta_g)$
 - where G is a differentiable generator function represented by a multilayer perceptron with parameters θ_g

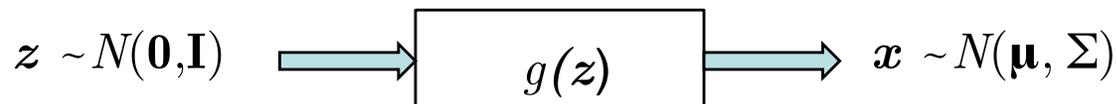


$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

$G(\mathbf{z})$ maps \mathbf{z} to \mathbf{x} which is distributed according to $p_g = p(\mathbf{x}|\mathbf{z})$

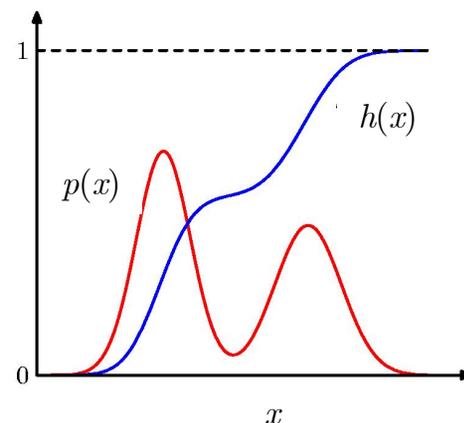
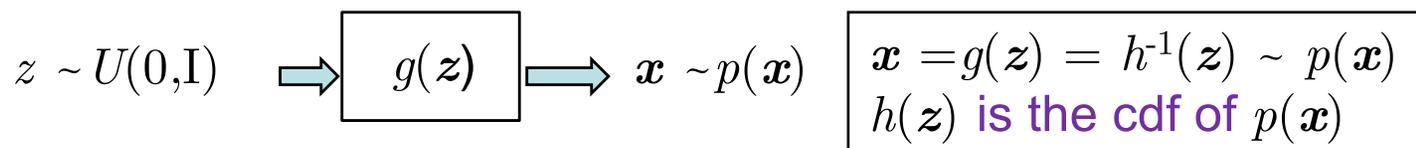
Differentiable Generator Functions

1. To sample from $\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$: $\mathbf{x} = g(\mathbf{z}) = \boldsymbol{\mu} + L\mathbf{z}$



- where L is the Cholesky decomposition of Σ
 - Decomposition: $\Sigma = LL^T$ where L is lower triangular

2. To sample from any $p(\mathbf{x})$



– Distribution is governed by

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(g^{-1}(\mathbf{x}))}{\left| \det\left(\frac{\partial g}{\partial \mathbf{z}}\right) \right|}$$

Denominator is determinant of Jacobian

• Alternatively, define $p(\mathbf{x}|\mathbf{z})$ and marginalize \mathbf{z} :

$$p(\mathbf{x}) = \mathbb{E}_{\mathbf{z}} p(\mathbf{x} | \mathbf{z})$$

e.g., by using a VAE

Generator and Discriminator Networks

1. Generator Network produces samples

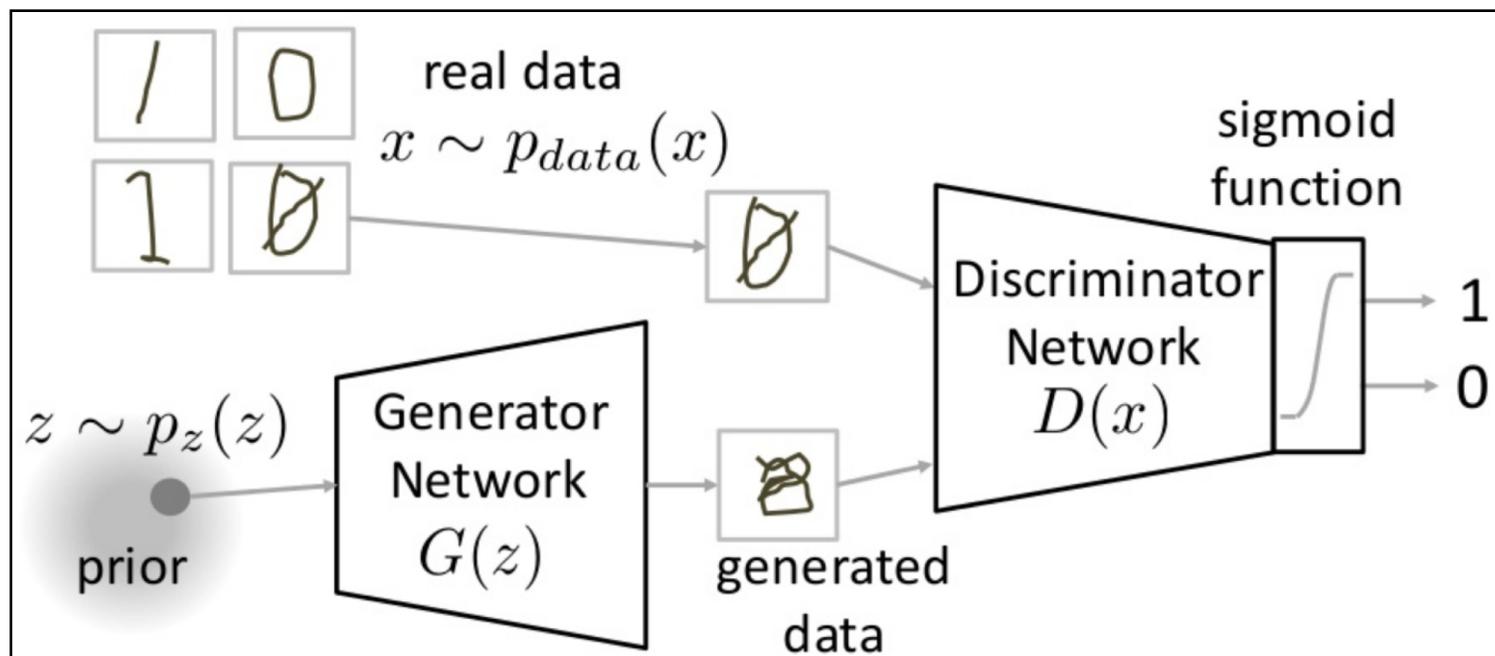
$$\boldsymbol{x} = G(\boldsymbol{z}; \theta^{(g)})$$

2. Discriminator Network

1. Attempts to distinguish samples drawn from training data and samples drawn from generator
2. Emits a probability value
 - $D(\boldsymbol{x}; \theta^{(d)})$ indicating the probability that \boldsymbol{x} is a real training example rather than a fake sample drawn from the model

Adversarial Modeling Framework

- We also define a second multilayer perceptron $D(\mathbf{x}; \theta_d)$ that outputs a single scalar
 - $D(\mathbf{x})$ represents the probability that \mathbf{x} came from the data rather than p_g



Training of D and G

1. We train D to maximize the probability of assigning the correct label to both training examples and samples from G
 2. We simultaneously train G to minimize $\log(1 - D(G(\mathbf{z})))$
- In other words, D and G play a two-player minimax game with *Value Function* $V(G, D)$

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

Minimax value explanation

1. D is accurate over real data by maximizing

$$\mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)]$$

– Given sample $G(z)$, $z \sim p_z(z)$, $D(G(z))$ is close to zero by maximizing $\mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$

2. G is trained to increase D 's probability for a fake example,

– i.e. minimize $\mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$

• Combined together D and G are playing a minimax game in which we optimize:

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

Optimal value of D

- To find best value for D consider

$$L(G, D) = \int_x (p_{\text{data}}(x) \log D(x) + p_g(x) \log(1 - D(x))) dx$$

- Since we want value of $D(x)$ to maximize $L(G, D)$, label $\tilde{x} = D(x)$, $A = p_{\text{data}}(x)$, $B = p_g(x)$

- What is inside integral is: \rightarrow

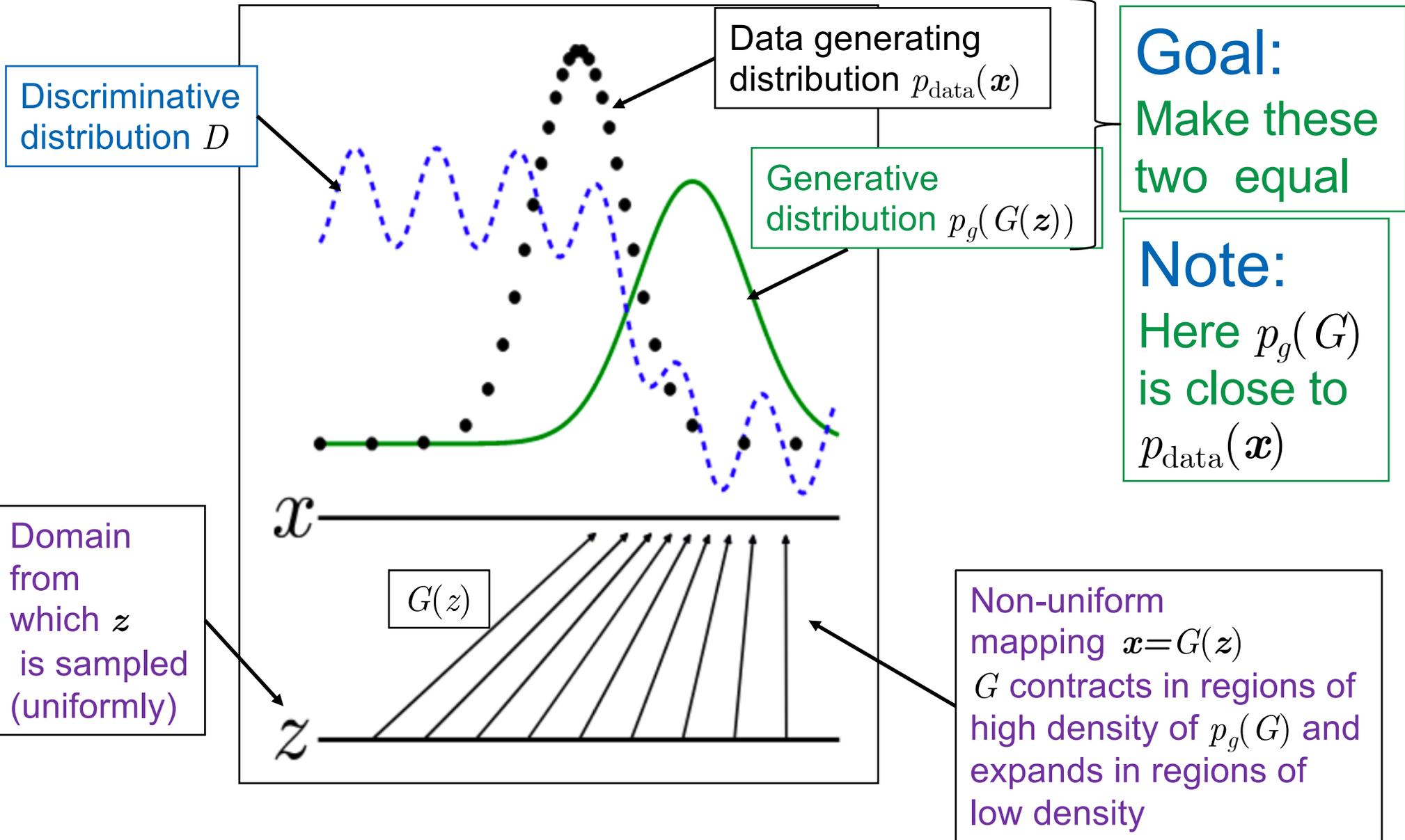
$$\begin{aligned} f(\tilde{x}) &= A \log \tilde{x} + B \log(1 - \tilde{x}) \\ \frac{df(\tilde{x})}{d\tilde{x}} &= A \frac{1}{\ln 10} \frac{1}{\tilde{x}} - B \frac{1}{\ln 10} \frac{1}{1 - \tilde{x}} \\ &= \frac{1}{\ln 10} \left(\frac{A}{\tilde{x}} - \frac{B}{1 - \tilde{x}} \right) \\ &= \frac{1}{\ln 10} \frac{A - (A + B)\tilde{x}}{\tilde{x}(1 - \tilde{x})} \end{aligned}$$

- Setting the derivative to zero:

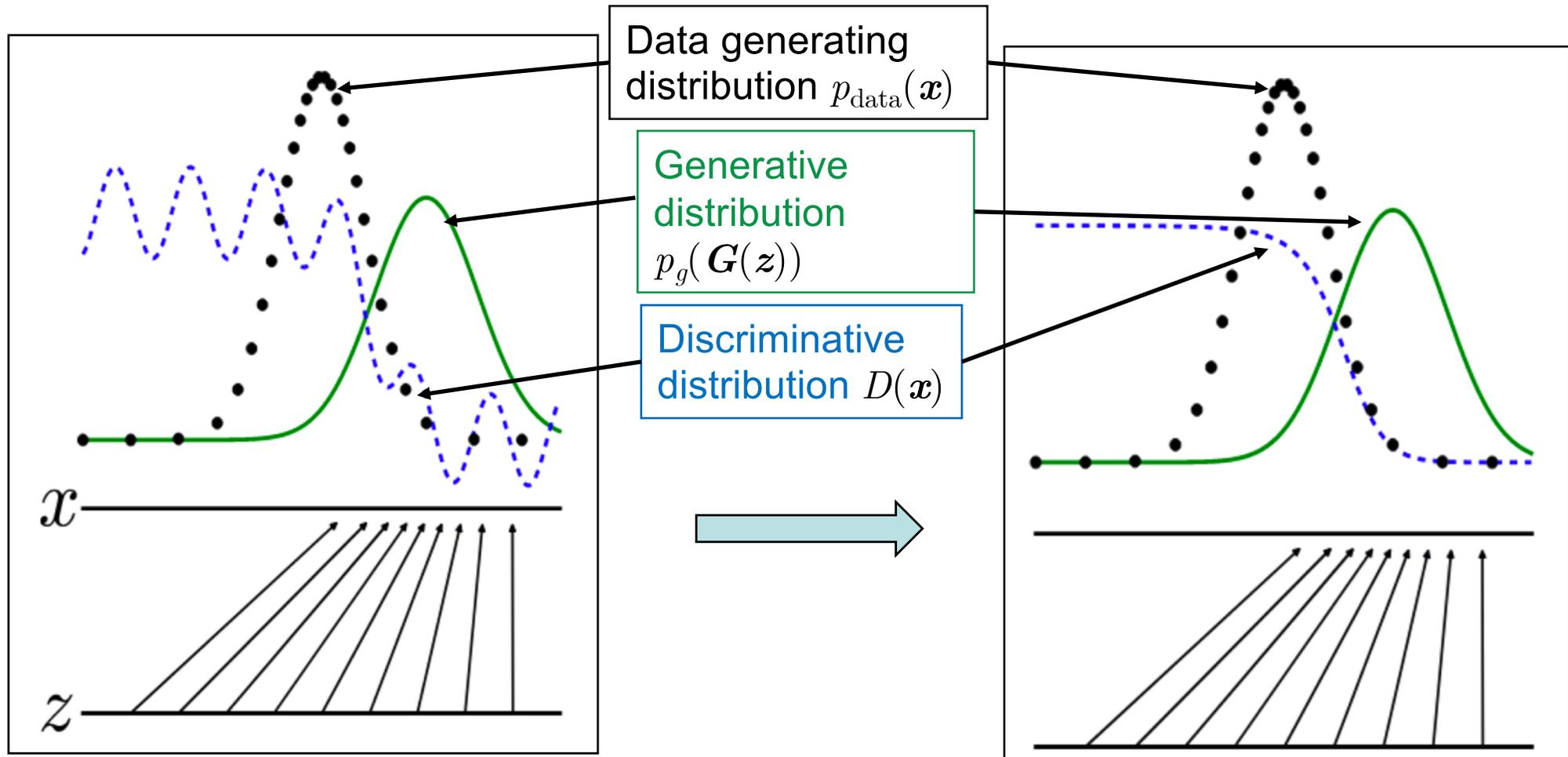
$$D^*(x) = \tilde{x}^* = \frac{A}{A+B} = \frac{p_r(x)}{p_r(x) + p_g(x)} \in [0, 1]$$

- Once G is trained to optimal, p_g is very close to p_{data}
- When $p_g = p_{\text{data}}$, $D^*(x)$ becomes $1/2$

Intuition in training of GAN



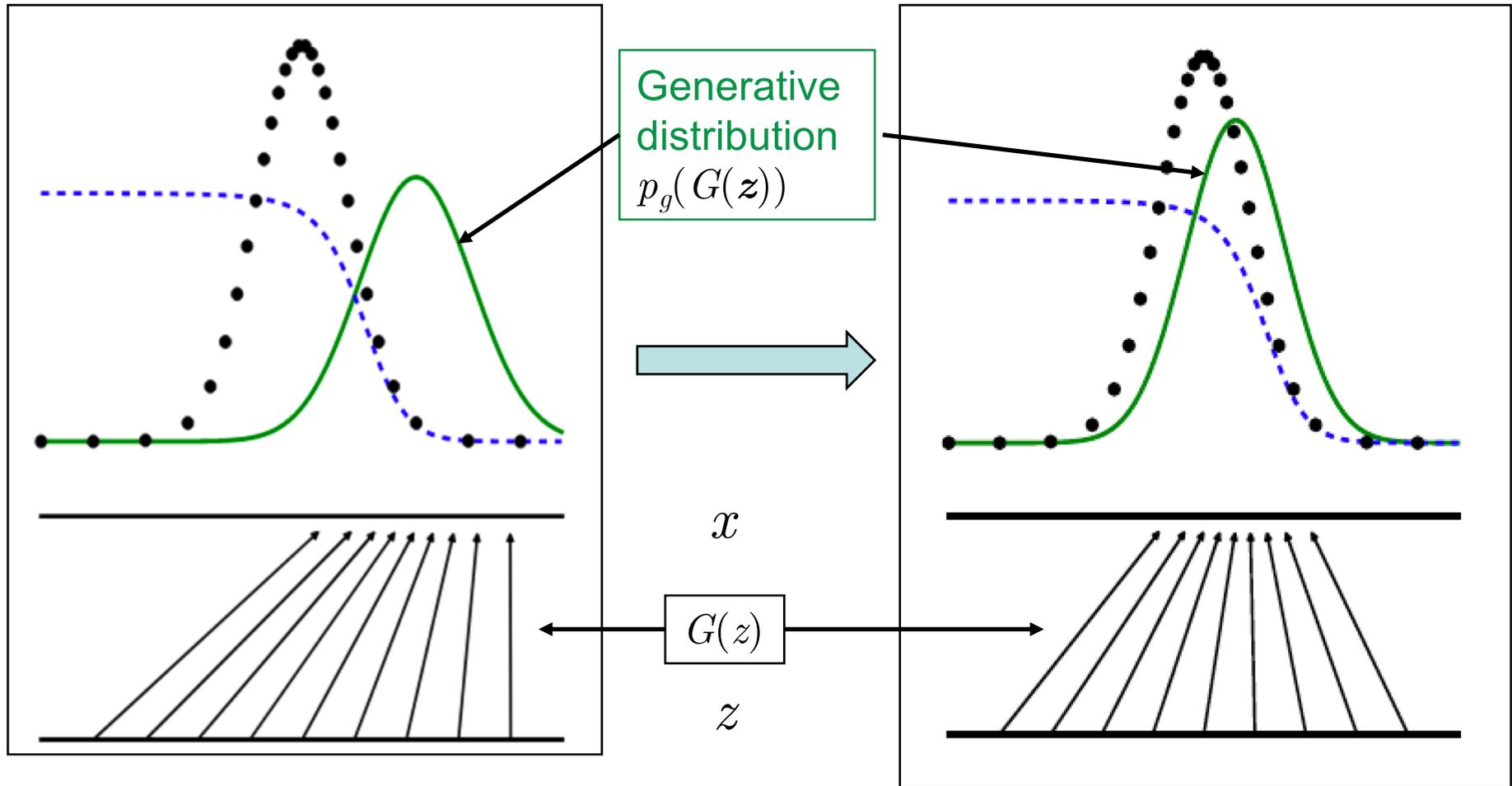
Update of D



D is trained to discriminate samples from data,

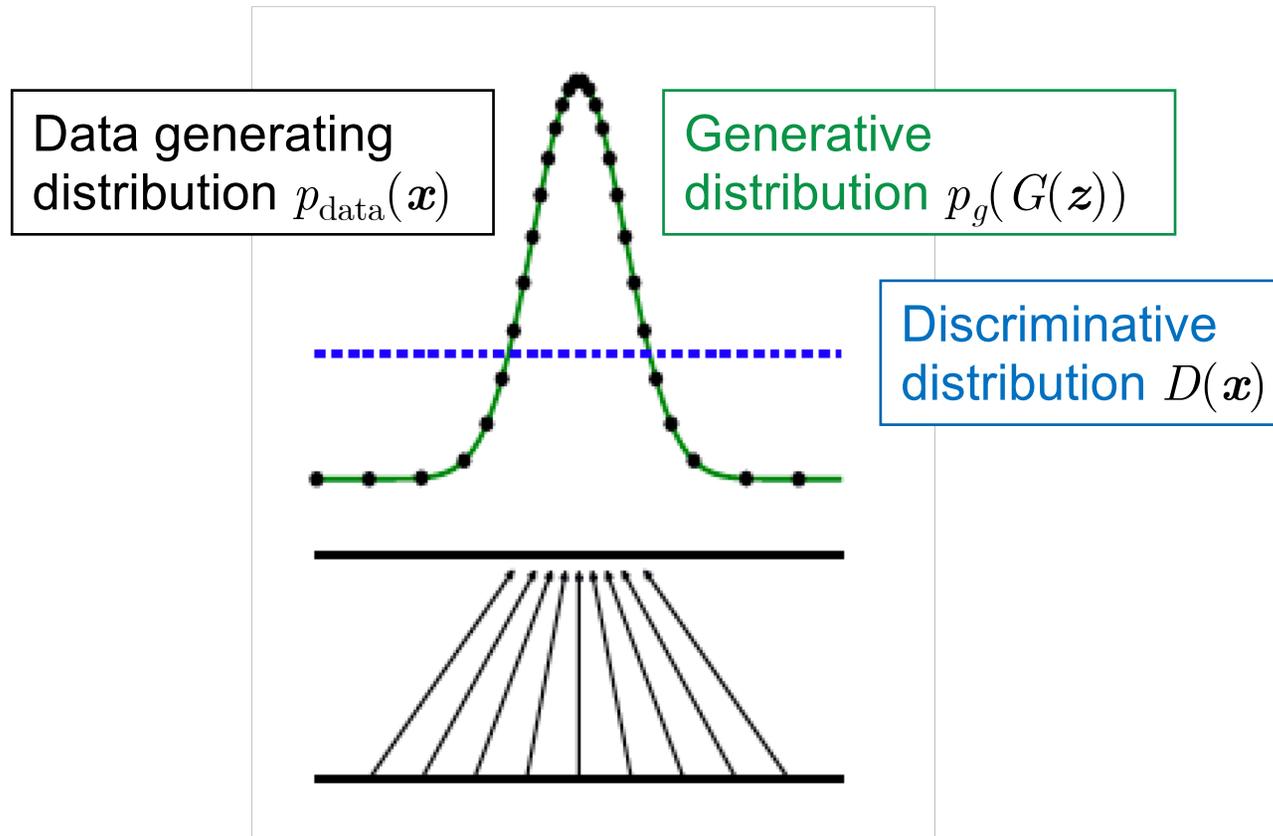
converging to
$$D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

Update of G



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

Final State



After several steps of training, if D and G 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., $D(\mathbf{x}) = 0.5$

GAN Objective

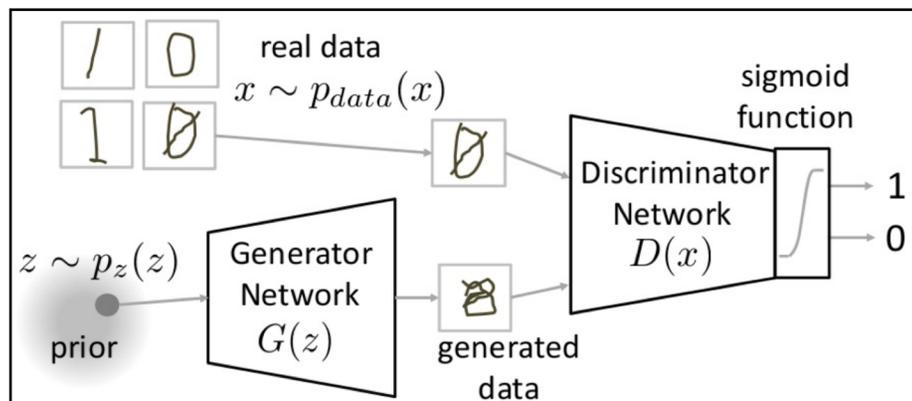
- Value Function of Minimax played by G, D

$$\min_G \max_D E_{x \sim p_X} [\log D(x)] + E_{z \sim p_Z} [\log(1 - D(G(z)))]$$

p_X : data distribution, represented by samples

$p_{G(z)}$: model distribution, z modeled as Gaussian

- We must implement the game using an iterative, numerical approach



Optimizing D to completion in inner loop of training is computationally prohibitive and results in overfitting in finite data sets.

Instead we alternate between k steps of optimizing D and one step of optimizing G

Interaction between G and D

- Weights are updated so as to
 - maximize the probability that any real data input x is classified as belonging to the real dataset, while
 - minimizing the probability that any fake image is classified as belonging to the real dataset
- In more technical terms,
 - Loss function maximizes the function $D(x)$, and it also minimizes $D(G(z))$

Alternating gradient updates

$$\min_G \max_D E_{x \sim p_X} [\log D(x)] + E_{z \sim p_Z} [\log(1 - D(G(z)))]$$

- Step 1: Fix G and perform a gradient step for maximization:

$$\max_D E_{x \sim p_X} [\log D(x)] + E_{z \sim p_Z} [\log(1 - D(G(z)))]$$

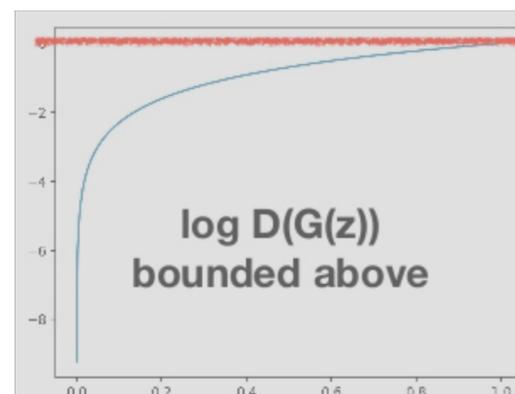
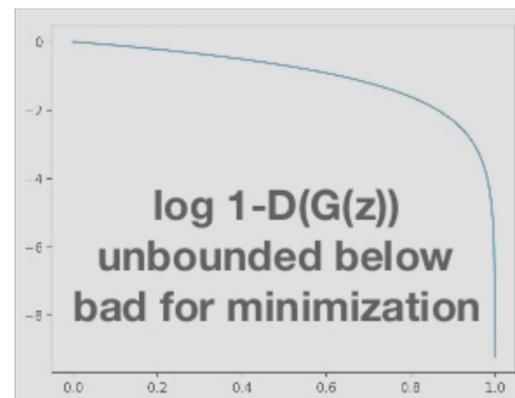
- Step 2: Fix D and perform a gradient step for minimization:

– In theory

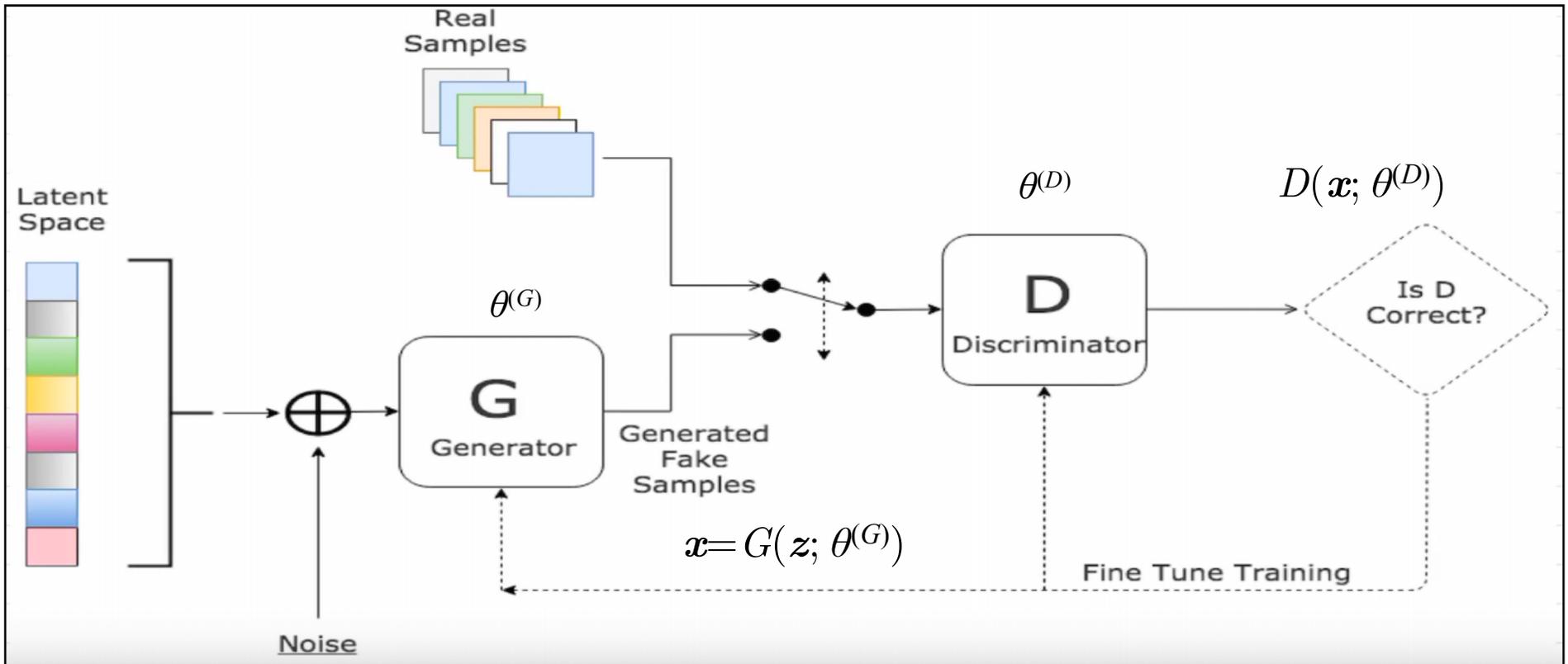
$$\min_G E_{z \sim p_Z} [\log(1 - D(G(z)))]$$

– In practice

$$\max_G E_{z \sim p_Z} [\log D(G(z))]$$



GAN Architecture



G and D are deep neural networks
Fine tune training consists of backpropagation

SGD for GAN

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:

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

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:

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

end for

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

Theoretical Results

- Theorem 1:
 - The global minimum of the virtual training criterion $C(G) = \max_D V(G, D)$ is achieved iff $p_g = p_{\text{data}}$
 - At that point $C(G)$ achieves the value $-\log 4$
- Convergence of Algorithm 1:
 - If G and D have enough capacity, and at each step of Algorithm 1, D is allowed to reach its optimum given G , and p_g is updated so to improve the criterion $\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$
 - Then p_g converges to p_{data}

GAN and Convexity

- Main motivation for GAN: No approximation
 - In inference or in partition function
- When $\max_d v(g, d)$ is convex in θ_g
 - (in the space of probability density functions)
 - the procedure is guaranteed to converge
- If g, d are neural nets and $\max_d v(g, d)$ nonconvex
 - Learning is difficult
 - From min-max theory, mode collapse is caused by the improper convergence to the optimum, due to the non-convexity and non-concavity of the objective function

GAN underfit: nonconvergence

- Simultaneous gradient descent on two player's costs is not guaranteed to reach an equilibrium
 - **Ex:** $v(a, b) = ab$
 - where Player 1 controls a and incurs cost ab , while Player 2 controls b and receives cost $-ab$
 - If each player makes small gradient steps, each player reducing own cost at expense of other,
 - then a and b go into a stable circular orbit, rather than arriving at the equilibrium point at the origin

- **Overfitting**
 - training loss \ll validation loss
- **Underfitting**
 - training loss \gg validation loss
- **Just right**
 - training loss \approx validation loss

Proposed Solutions

1. Game is no longer zero-sum
 - Same expected gradient as maximum likelihood learning
 - Does not work well in practice
2. G increases probability that D makes a mistake rather than aiming to decrease probability that D makes right prediction
 - Motivated by keeping derivatives large
- Stabilization of GAN learning an open problem

References

- <https://www.kdnuggets.com/2019/04/rise-generative-adversarial-networks.html?fbclid=IwAR3wmr2-MbiNEjoFGu6QoocPmOCmAuPYPBfmolDQcolJWThnXRiVP7iKGSk>
- <https://www.slideshare.net/mlreview/tutorial-on-theory-and-application-of-generative-adversarial-networks>