

Generative Adversarial Networks

: From Theory to Practice

Sunmook Choi

Dept. of Mathematics
Korea University

Generative Adversarial Networks

Given a dataset, we want to **generate new data** which seem to be obtained from the dataset, i.e., given p_{data} , we want to learn the **generator's distribution** p_g over data \mathbf{x} s.t. $p_{\text{data}} = p_g$.

How?

- Define a prior on input noise/latent variable $p_z(\mathbf{z})$, and
- represent a mapping to a data space as $G(\mathbf{z}; \theta_g)$, which is the generator.
- Define the discriminator $D(\mathbf{x}; \theta_d)$ that outputs the probability that \mathbf{x} came from the data distribution p_{data} rather than p_g .
- Train D to maximize the probability of assigning the correct label to both training examples (real, label=1) and samples from G (fake, label=0).
- Simultaneously train G to minimize $\log(1 - D(G(\mathbf{z})))$.

Objective Function

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

Recall: Binary Cross Entropy Loss for Logistic Regression

For binary dataset $\{\mathbf{x}^n, y^n\}_{n=1}^N$ where $y^n \in \{0, 1\}$ with hypothesis model h_θ ,

$$BCE(\theta) = \frac{1}{N} \sum_{n=1}^N \left[-y^n \log(h_\theta(\mathbf{x}^n)) - (1 - y^n) \log(1 - h_\theta(\mathbf{x}^n)) \right]$$

- For D , $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$ should have label 1, and $G(\mathbf{z}) \sim p_g(\mathbf{x})$ should have label 0. Hence, maximizing $V(D, G)$ w.r.t. D is equivalent to minimizing BCE.

Objective Function

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

Recall: Binary Cross Entropy Loss for Logistic Regression

For binary dataset $\{\mathbf{x}^n, y^n\}_{n=1}^N$ where $y^n \in \{0, 1\}$ with hypothesis model h_θ ,

$$BCE(\theta) = \frac{1}{N} \sum_{n=1}^N \left[-y^n \log(h_\theta(\mathbf{x}^n)) - (1 - y^n) \log(1 - h_\theta(\mathbf{x}^n)) \right]$$

- Practically, minimizing $V(D, G)$ w.r.t. G may not provide sufficient gradient for G early in learning. Instead, we train G to maximize $\log D(G(\mathbf{z}))$, i.e., **minimize $-\log D(G(\mathbf{z}))$** .
- Hence, minimizing $-\log D(G(\mathbf{z}))$ w.r.t. G is equivalent to minimizing BCE.

Training Algorithm

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:

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

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:

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

end for

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

Theoretical Results

Proposition 1. For G fixed, the optimal discriminator D is $D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$.

$$\because V(D, G) = \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{z}} p_z(\mathbf{z}) \log(1 - D(G(\mathbf{z}))) d\mathbf{z} = \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x}$$

For any $(a, b) \in \mathbb{R}^2 \setminus (0, 0)$, the function $y \mapsto a \log(y) + b \log(1 - y)$ achieves its maximum in $[0, 1]$ at $\frac{a}{a+b}$.

Using the result above, the objective function w.r.t. G can be reformulated as follows:

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

Theoretical Results

Theorem 1. The global minimum of the virtual training criterion $C(G)$ is achieved if and only if $p_g = p_{\text{data}}$.

\because Recall that $D_{KL}(p \parallel q) = \mathbb{E}_{x \sim p} \left[\log \frac{p(x)}{q(x)} \right]$ and $D_{JS}(p, q) = \frac{1}{2} D_{KL} \left(p \parallel \frac{p+q}{2} \right) + \frac{1}{2} D_{KL} \left(q \parallel \frac{p+q}{2} \right)$.

$$\begin{aligned} C(G) &= \mathbb{E}_{x \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{x \sim p_g} \left[\log \frac{p_g(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right] \\ &= D_{KL} \left(p_{\text{data}} \parallel \frac{p_{\text{data}} + p_g}{2} \right) + D_{KL} \left(p_g \parallel \frac{p_{\text{data}} + p_g}{2} \right) - 2 \log 2 \\ &= D_{JS}(p_{\text{data}}, p_g) - \log 4 \end{aligned}$$

Since the Jensen-Shannon divergence is a metric, we have shown that the global minimum of $C(G)$ is $-\log 4$ and that the only solution is $p_g = p_{\text{data}}$, i.e., **the generative model perfectly replicates the data generating process.**