

Generative Adversarial Networks

Setup: Given $D := \{x_1, \dots, x_n\}$, $x_i \in \mathbb{R}^d$, $x_i \sim q(x)$ empirical data distribution

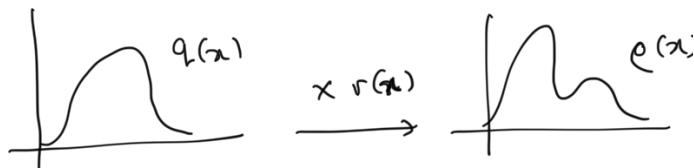
$$p_\theta(x) = \int p(x, z) dz = \int \underbrace{p(x|z)}_{\mathcal{N}(x | \mu_\theta(z), \Sigma_\theta(z))} \underbrace{p(z)}_{\mathcal{U}(0, 1)} dz, \quad p(z|x) \approx q_\phi(z)$$

$$\underbrace{L(\theta, \phi)}_{\text{MLE}} = \underbrace{-\log p_\theta(x)}_{\text{MLE}} \stackrel{\text{ELBO}}{\leq} \text{KL}[q_\phi(z|x) \| p(z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)]$$

$$\theta^*, \phi^* = \underset{\theta, \phi}{\text{argmin}} \text{KL}[q(x) \| p_\theta(x)] = \mathbb{E}_{x \sim q(x)} \left[\log \frac{q(x)}{p_\theta(x)} \right]$$

Density-ratio estimation:

$$r(x) = \frac{p(x)}{q(x)}$$



$$x_i \sim p(x)$$

$$x_j \sim q(x)$$

$$X = \begin{bmatrix} \overbrace{x_1, \dots, x_n}^{p(x)} & \overbrace{x'_1, \dots, x'_n}^{q(x)} \end{bmatrix} \quad \begin{matrix} 2n \times d \\ 2n \times 1 \end{matrix}$$

$$y = \begin{bmatrix} +1, \dots, +1 & -1, \dots, -1 \end{bmatrix}$$

$$p(x) = p(x | y = +1)$$

$$q(x) = p(x | y = -1)$$

$$p(x|y) = \frac{p(y|x) p(x)}{p(y)}$$

$$r(x) = \frac{p(x)}{q(x)} = \frac{p(x|y=+1)}{p(x|y=-1)} \stackrel{\text{Bayes}}{=} \frac{\frac{p(y=+1|x)p(x)}{p(y=+1)}}{\frac{p(y=-1|x)p(x)}{p(y=-1)}}$$

$$= \frac{\cancel{p(y=-1)}}{\cancel{p(y=+1)}} \cdot \frac{p(y=+1|x)}{p(y=-1|x)} = \frac{p(y=+1|x)}{1 - p(y=+1|x)} \approx \frac{S_+(x)}{1 - S_+(x)}$$

GANs:

Goal is to construct a model for $p_\theta(x)$.

$$G_\theta(z) = x, \quad \overset{\sim \mathcal{N}(0, I)}{z} \sim p(z), \quad z \in \mathbb{R}^q, \quad x \in \mathbb{R}^d$$

$$G_\theta: \mathbb{R}^q \rightarrow \mathbb{R}^d, \quad x \sim p_\theta(x)$$

We want to train this generative mapping such that the distribution $p_\theta(x)$ is as close as possible to $q(x)$ (empirical dist. of the observed data)

Unlike MLE approaches, here we introduce a discriminator

$$D_\phi: \mathbb{R}^d \rightarrow [0, 1], \quad D_\phi(x)$$

Training:

$$\theta^*, \phi^* = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim q(x)} [\log D_\phi(x)] + \mathbb{E}_{z \sim p(z)} [1 - \log D_\phi(\overset{\text{true}}{\uparrow} G_\theta(z))] \quad \text{fals}$$

Algorithm:

for t iter

for (12) steps do :

sample $z_1, \dots, z_m \sim p(z)$

sample $x_1, \dots, x_m \sim q(x)$

$$\phi_{n+1} = \phi_n + \eta \nabla_{\phi} \left\{ \frac{1}{m} \sum_{i=1}^m \log(D_{\phi}(x_i)) + \log[1 - D_{\phi}(G_{\theta}(z_i))] \right\}$$

while θ is kept fixed.

binary cross-entropy

keep θ
params
fixed
 \uparrow
 η

for (t) steps do :

sample $z_1, \dots, z_m \sim p(z)$

$$\theta_{n+1} = \theta_n - \eta \nabla_{\theta} \left\{ \frac{1}{m} \sum_{i=1}^m \log[1 - D_{\phi}(G_{\theta}(z_i))] \right\}$$

while ϕ is kept fixed.

Remarks :

1.) For $G_{\theta}(z)$ fixed, the optimal discriminator is :

$$D_{\phi}^*(x) = \frac{q(x)}{q(x) + p_{\theta}(x)} \quad \checkmark$$

$$J(\theta) = \max_{\phi} \mathcal{L}(\theta, \phi) = \mathbb{E}_{x \sim q(x)} \left[\log \frac{q(x)}{q(x) + p_{\theta}(x)} \right] + \mathbb{E}_{x \sim p_{\theta}(x)} \left[\log \frac{p_{\theta}(x)}{q(x) + p_{\theta}(x)} \right]$$

The global minimum of $J(\theta)$ can be achieved iff $p_{\theta}(x) =$

$$J(\theta^*) = -\log 4$$

Goodfellow et al 20

$J(\theta)$ can be re-written as :

$$J(\theta) = -\log 4 + \underbrace{\text{KL}[q(x) \parallel \frac{q(x) + p_{\theta}(x)}{2}] + \text{KL}[p_{\theta}(x) \parallel \frac{q(x) + p_{\theta}(x)}{2}]}$$

Jensen-Shannon entropy

$$J(\theta) = -\log 4 + 2 \text{JSD}[q(x) \| p_\theta(x)]$$

2.) If $G_\theta(z)$ and $D_\phi(x)$ have enough capacity, and at each step of the training algorithm $D_\phi(x)$ is allowed to reach its optimum, and $p_\theta(x)$ is updated to improve the following criterion:

$$\underbrace{\mathbb{E}_{x \sim q(x)} [\log D_\phi^*(x)] + \mathbb{E}_{z \sim p(z)} [\log (1 - D_\phi^*(G_\theta(z)))]}$$

then $p_\theta(x) \rightarrow q(x)$.