

# CS 446: Machine Learning

## Homework 11

Due on Tuesday, April 24, 2018, 11:59 a.m. Central Time

### 1. [8 points] Generative Adversarial Network (GAN)

- (a) What is the cost function for classical GANs? Use  $D_w(x)$  as the discriminator and  $G_\theta(x)$  as the generator.

Your answer:  $V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D_w(x)] + \mathbb{E}_z [\log(1 - D_w(G_\theta(z)))]$

- (b) Assume arbitrary capacity for both discriminator and generator. In this case we refer to the discriminator using  $D(x)$ , and denote the distribution on the data domain induced by the generator via  $p_G(x)$ . State an equivalent problem to the one asked for in part (a), by using  $p_G(x)$  and the ground truth data distribution  $p_{\text{data}}(x)$ .

Your answer:

$$V(D, G) = \int_x p_{\text{data}}(x) \log D(x) \, dx + \int_z p_z(z) \log(1 - D(G(z))) \, dz$$

$$x = G(z) \implies z = G^{-1}(x) \implies dz = [G^{-1}]'(x) dx$$

$$V(D, G) = \int_x p_{\text{data}}(x) \log D(x) \, dx + \int_x p_z(G^{-1}(x)) [G^{-1}]'(x) \log(1 - D(x)) \, dx$$

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

- (c) Assuming arbitrary capacity, derive the optimal discriminator  $D^*(x)$  in terms of  $p_{data}(x)$  and  $p_G(x)$ .

You may need the Euler-Lagrange equation:

$$\frac{\partial L(x, D, \dot{D})}{\partial D} - \frac{d}{dx} \frac{\partial L(x, D, \dot{D})}{\partial \dot{D}} = 0$$

where  $\dot{D} = \partial D / \partial x$ .

Your answer:

$$\max_D V(D, G) = \max_D \int_x p_{data}(x) \log D(x) + p_G(x) \log(1 - D(x)) dx$$

Applying Euler-Lagrange equation to find the optimal discriminator  $D^*(x)$

let  $L(x, D) = p_{data}(x) \log D(x) + p_G(x) \log(1 - D(x))$

$$\frac{\partial L(x, D)}{\partial D} - \frac{d}{dx} \frac{\partial L(x, D)}{\partial \dot{D}} = 0$$

$$\Rightarrow \frac{\partial}{\partial D(x)} (p_{data}(x) \log D(x) + p_G(x) \log(1 - D(x))) = 0$$

$$\frac{p_{data}(x)}{D(x)} - \frac{p_G(x)}{1 - D(x)} = 0$$

$$D^*(x) = \frac{p_{data}(x)}{p_{data} + p_G(x)}$$

- (d) Assume arbitrary capacity and an optimal discriminator  $D^*(x)$ , show that the optimal generator,  $G^*(x)$ , generates the distribution  $p_G^* = p_{data}$ , where  $p_{data}(x)$  is the data distribution

You may need the Jensen-Shannon divergence:

$$\text{JSD}(p_{data}, p_G) = \frac{1}{2} D_{KL}(p_{data}, M) + \frac{1}{2} D_{KL}(p_G, M) \quad \text{with} \quad M = \frac{1}{2}(p_{data} + p_G)$$

Your answer:

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

$$C(G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D^*(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D^*(G(z)))]$$

$$C(G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D^*(x)] + \mathbb{E}_{x \sim p_G} [\log(1 - D^*(x))]$$

$$C(G) = \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \left( \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G^*(x)} \right) \right] + \mathbb{E}_{x \sim p_g} \left[ \log \left( \frac{p_G^*(x)}{p_{\text{data}}(x) + p_G^*(x)} \right) \right]$$

$$C(G) = \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \left( \frac{p_{\text{data}}(x)}{\frac{p_{\text{data}}(x) + p_G^*(x)}{2}} \right) \right] + \mathbb{E}_{x \sim p_g} \left[ \log \left( \frac{p_G^*(x)}{\frac{p_{\text{data}}(x) + p_G^*(x)}{2}} \right) \right] - \log(4)$$

$$C(G) = D_{KL}(p_{\text{data}}, M) + D_{KL}(p_G^*, M) - \log(4)$$

$$C(G) = 2 \times \text{JSD}(p_{\text{data}}, p_G^*) - \log(4)$$

Since, the Jensen - Shannon divergence between two distributions are non negative, and zero iff they are equal. Hence, we have shown that  $C^* = -\log(4)$  is the global minimum of the of  $C(G)$  and  $p_G^* = p_{\text{data}}$

Alternatively:

We can optimize it point-wise, let  $p_G^* = a$  and  $p_{\text{data}} = x$  and taking derivative with respect to x. we will get :

$$\frac{\partial}{\partial x} \left( a \log \left( \frac{2a}{a+x} \right) + x \log \left( \frac{2x}{a+x} \right) \right) = \log(2x) - \log(a+x) = 0$$

Since the above is monotonic, hence there will be only one solution i.e.  $a = x$  i.e.  $p_G^* = p_{\text{data}}$

- (e) More recently, researchers have proposed to use the Wasserstein distance instead of divergences to train the models since the KL divergence often fails to give meaningful information for training. Consider three distributions,  $\mathbb{P}_1 \sim U[0, 1]$ ,  $\mathbb{P}_2 \sim U[0.5, 1.5]$ , and  $\mathbb{P}_3 \sim U[1, 2]$ . Calculate  $D_{KL}(\mathbb{P}_1, \mathbb{P}_2)$ ,  $D_{KL}(\mathbb{P}_1, \mathbb{P}_3)$ ,  $\mathbb{W}_1(\mathbb{P}_1, \mathbb{P}_2)$ , and  $\mathbb{W}_1(\mathbb{P}_1, \mathbb{P}_3)$ , where  $\mathbb{W}_1$  is the Wasserstein-1 distance between distributions.

Your answer: KL Divergence is defined as:

$$D_{KL}(U, P) = \int_x u(x) \log \frac{u(x)}{p(x)}$$

Therefore,

$$D_{KL}(\mathbb{P}_1, \mathbb{P}_2) = D_{KL}(\mathbb{P}_1, \mathbb{P}_3) = \infty$$

Definition of Wassertein-1 distance:

$$W(p_r, p_\theta) = \inf_{\gamma \in \pi} \iint_{x, y} \|x - y\| \gamma(x, y) \, dx \, dy = \inf_{\gamma \in \pi} \mathbb{E}_{x, y \sim \gamma} [\|x - y\|].$$

Where  $\pi$  is the joint probability of  $x$  and  $y$  over all possible paths  $\gamma$ . Therefore,

$$\mathbb{W}_1(\mathbb{P}_1, \mathbb{P}_2) = 0.5 \text{ and } \mathbb{W}_1(\mathbb{P}_1, \mathbb{P}_3) = 1$$