

An alternative view of GAN

Consider the following process :

① the machine flips a (possibly bent) coin

② depending on the outcome :

1) if heads, the machine shows a real image

2) if tails, the machine shows a generated image

This process is repeated several times.

You're asked to guess the probability that the coin flip outcome is a tail.

What would the generative model do to fool you on this task?

This process is equivalent to the following augmented generative model

$$S \sim \hat{P}(S) = \pi^S (1-\pi)^{1-S}$$

$$x|S \sim \hat{P}(x|S) = \begin{cases} p(x), & S=1 \\ p_0(x), & S=0 \end{cases}$$

We train this model by

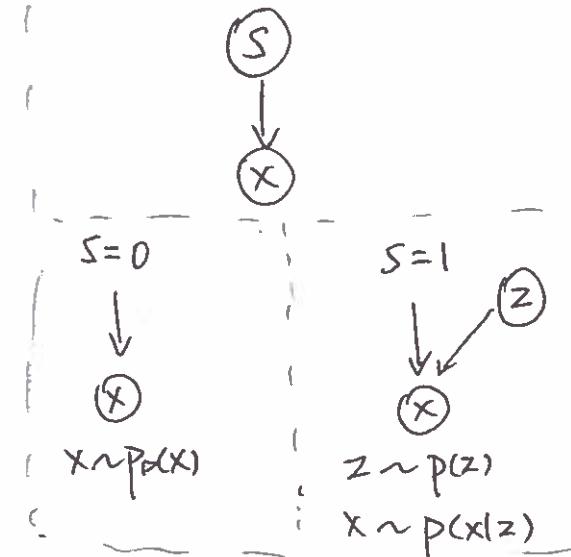
$$\min_{p(x)} I[S; x] = KL[\hat{P}(S, x) || \hat{P}(S)\hat{P}(x)]$$

when $I[S; x] = 0$:

observation x tells us nothing about the coin flip outcome S (So we can only do random guess).

Notes: ① $I[S; x] = 0 \Leftrightarrow p(x) = p_0(x)$

② when $\pi = \frac{1}{2}$, $I[S; x] = JS[p(x) || p_0(x)]$



$$\min_{p(x)} I[s; x] \quad (\text{intractable})$$

Bringing in the idea of variational approximation:
 (similar to variational information maximization)

$$\begin{aligned} I[s; x] &\geq I[s; x] - \mathbb{E}_{\hat{p}(x)} [\text{KL}[\hat{p}(s|x) \| q(s|x)]] \\ &= H[s] - H[s|x] + H[s|x] + \mathbb{E}_{\hat{p}(x,s)} [\log q(s|x)] \\ &= H[s] + \mathbb{E}_{\hat{p}(x,s)} [\log q(s|x)] \end{aligned}$$

Now the objective becomes

$$\min_{p(x)} \max_{q(s|x)} H[s] + \mathbb{E}_{\hat{p}(x,s)} [\log q(s|x)] \quad (\text{tractable})$$

$$\begin{aligned} \mathbb{E}_{\hat{p}(x,s)} [\log q(s|x)] &= \underbrace{\pi (1-\pi) \mathbb{E}_{p_0(x)} [\log q(s=0|x)] + \pi \mathbb{E}_{p(x)} [\log (1-q(s=0|x))]}_{\text{GAN's objective (when } \pi = \frac{1}{2} \text{)}} \\ &= \text{GAN's objective (when } \pi = \frac{1}{2} \text{)} \end{aligned}$$

ALI as an augmented generative model

$$S \sim p(s) = \pi^s (1-\pi)^{1-s}$$

$$x, z | s \sim \hat{p}(x, z | s) = \begin{cases} p(z)p(x|z), & s=1 \\ p_0(x)q(z|x), & s=0 \end{cases}$$

we train this model by

$$\min_{p(x|z), q(z|x)} I[S; \{x, z\}]$$

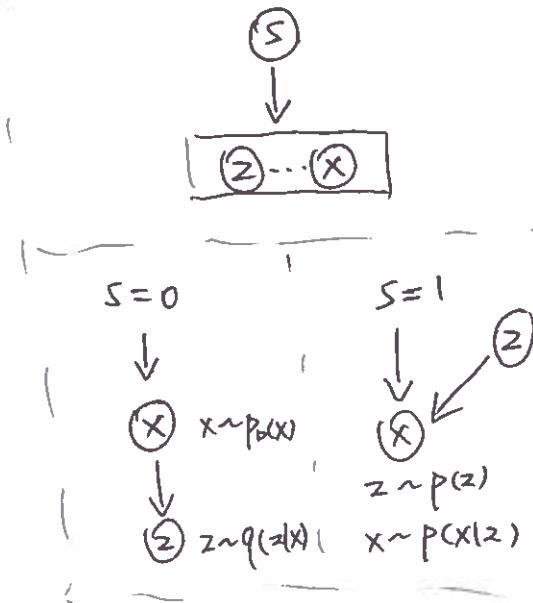
Notes: under this model:

$$\textcircled{1} \quad I[S; \{x, z\}] = I[S, x] + I[S, z|x]$$

$$I[S, z|x] = \mathbb{E}_{p(x)} [\text{KL}[\hat{p}(s, z|x) || \hat{p}(s|x)\hat{p}(z|x)]]$$

$$\text{so } I[S, z|x] = 0 \iff p(z|x) = q(z|x)$$

so min. over $q(z|x)$ make sense!



To see why $I[S, z|x] = 0 \Leftrightarrow p(z|x) = q(z|x)$:

$$I[S, z|x] = 0 \Leftrightarrow \hat{p}(s, z|x) = \hat{p}(s|x)\hat{p}(z|x)$$

$$\hat{p}(s, z|x) = \hat{p}(s|x)\hat{p}(z|s,x) = \begin{cases} q(z|x), & s=0 \\ p(z|x), & s=1 \end{cases}$$

$$\hat{p}(z|x) = \hat{p}(s=0|x)q(z|x) + \hat{p}(s=1|x)p(z|x)$$

$$\hat{p}(z|x) = \hat{p}(z|s,x), s=0,1 \Leftrightarrow p(z|x) = q(z|x)$$

② ALI's objective function:

$$\min_{p(x,z), q(z|x)} \max_{q(s|x,z)} I[S; \{x, z\}] - \mathbb{E}_{\hat{p}(x,z)} [\text{KL}[\hat{p}(s|x,z) \| q(s|x,z)]]$$
$$= H[S] + \mathbb{E}_{\hat{p}(s,x,z)} [\log q(s|x,z)]$$