

Using LVM for generation:

→ sample $z \sim p(z)$

→ sample $x \sim p(x|z)$

"similar at the population level" → becomes hard to distinguish real/fake

GAN is a game. Train a network to "guess" real/fake.
2-player. "discriminator": serves as loss function for "generator".

Objective of generator $G(z)$: make discriminator $D(x) = 0.5$ for all generated x (cannot tell).

↳ generate realistic images

↳ generate all possible realistic images (i.e. learn distribution).

zero-sum

classic GAN 2-player game:

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

$$\approx \frac{1}{N} \sum_{i=1}^N \log D(x_i) \quad x_i \in \mathcal{D}_r \quad \approx \frac{1}{N} \sum_{j=1}^N \log (1 - D(x_j)) \quad x_j = G(z_j)$$

more clearly: $\min_{\theta} \max_{\phi} V(\theta, \phi) = \mathbb{E}_{x \sim p_{data}(x)} [\log D_{\phi}(x)] + \mathbb{E}_{z \sim p(z)} [\log (1 - D_{\phi}(G_{\theta}(z)))]$

SGD → $\theta \leftarrow \theta + \lambda \nabla_{\theta} V(\theta, \phi) \approx \nabla_{\theta} \left(\frac{1}{N} \sum_{i=1}^N \log D_{\phi}(x_i) + \frac{1}{N} \sum_{j=1}^N \log (1 - D_{\phi}(x_j)) \right) \leftarrow \text{CE loss}$

$\phi \leftarrow \phi - \lambda \nabla_{\phi} V(\theta, \phi)$

$\approx \nabla_{\phi} \left(\frac{1}{N} \sum_{i=1}^N \log (1 - D_{\phi}(G_{\theta}(z_j))) \right)$

optimal discriminator (Bayes classifier) = $D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_G(x)} \quad x = G(z), \quad z \sim p(z)$

objective for G :

$V(D_G^*, G) =$

$\mathbb{E}_{p_{data}(x)} [\log p_{data}(x)] - \log(p_{data}(x) + p_G(x)) +$

$\mathbb{E}_{p_G(x)} [\log p_G(x)] - \log(p_{data}(x) + p_G(x))]$

$= D_{JS}(p_{data} || p_G)$

practical generator loss: $\mathbb{E}_{z \sim p(z)} [-\log D_{\phi}(G_{\theta}(z))]$

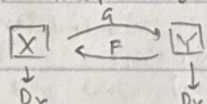
"maximize probability that image is real"

↳ better behaving gradients.

uses of GANs: conditional GAN (labellings)

CycleGAN ("translating")

↳ two (conditional) generators and two corresponding discriminators.



"Earth mover's distance"

↳ incorporates cycle-consistency loss

Given $x \rightarrow \hat{y} \rightarrow \hat{x}$, how close are x/\hat{x} ?

WGAN: accounts for "distance" b/w p_{data} and p_G .

$W(p_{data}, p_G) = \inf_{\gamma} \mathbb{E}_{x \sim p_{data}, y \sim \gamma(x, y)} [\|x - y\|]$

$\gamma_X(x) = p_{data}(x)$

$\gamma_Y(y) = p_G(y)$

really complex...

Kantorovich-Rubinstein duality

$= \sup_{f: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E}_{p_{data}} [f(x)] - \mathbb{E}_{p_G(x)} [f(x)]$

$\|f\|_L \leq 1$

set of all 1-Lipschitz scalar functions

$|f(x) - f(y)| \leq \|x - y\|$

(bounded slope, not too steep).

i.e. easier to train!

improved GAN techniques:

→ least-squares GAN (LSGAN)

↳ discriminator outputs real-valued num

→ Wasserstein GAN (WGAN)

↳ discriminator is Lipschitz-continuous

→ Gradient penalty

↳ discriminator is constrained to be continuous even harder

→ spectral norm

↳ discriminator is really constrained to be continuous

→ instance noise

↳ try to get $p_{data}(x)/p_G(x)$ overlap.

good choices today!

Gradient penalty: bounded slope.

update θ using gradient of

$\mathbb{E}_{x \sim p_{data}} [f_{\theta}(x)] - \frac{1}{2} (\|\nabla_x f_{\theta}(x)\|_2 - 1)^2 - \mathbb{E}_{z \sim p(z)} [f_{\theta}(G(z))]$

make norm of gradient close to 1.

$\sigma(W) = \max_{h \neq 0} \frac{\|Wh\|}{\|h\|}$

$= \max_{\|h\| \leq 1} \|Wh\|$

largest of W .

spectral norm: bound the Lipschitz constant in terms of singular values of each W_i .

↳ max slope of $Wx+b$ is spectral norm: $W_i \leftarrow W_i / \sigma(W_i)$.