

Towards Principled Methods for Training GANs

1. Discriminator reaches zero loss, ruins generator $\ddot{=}$
2. Ruins generator gradients $\ddot{=}$
3. Adding continuous noise helps $\ddot{=}$

0. Preliminaries

- g_0 generator, D discriminator, $g_0: Z \rightarrow X$
- P_g vs $P_r \leftarrow$ real dist. $z \sim p(z)$
- $L(D, g_0) = E_{P_r}[\log D(x)] + E_{P_g}[\log(1 - D(x))]$ ~~JSD~~ JS D
- Only need 2nd term for training g_0

1. Discriminator Reaches 0 loss

Key - P_r and P_g are on low dim manifolds
 - telling manifolds apart is easy

Lemma: $f: Z \rightarrow X$, f is "nice", $f(Z)$ is countable union of manifolds with dimension $\leq \dim(Z)$

PF: not obvious. Peano curves are counterexample

$f =$ Affine maps + Diffeomorphisms \checkmark

Lemma: w.h.p 2 manifolds intersect in lower dimensional manifolds

Thm: perfect discriminator: force it to be smooth on the manifold X

$$D^*(x) = 0 \text{ or } 1$$

$$\nabla_x D^*(x) = 0$$

at all input places

2. Sad Generator

- Investigate $\nabla_{\theta} E_{z \sim p(z)} [\log(1 - D(g_{\theta}(z)))]$

$$\nabla_{\theta} E_z [\log(1 - D(g_{\theta}(z)))] \leq E_z \left[\frac{\|\nabla_{\theta} D(g_{\theta}(z))\|^2}{|1 - D(g_{\theta}(z))|^2} \right]$$

$$\leq E_z \left[\frac{\|\nabla_x D(g_{\theta}(z))\|^2 \|J_{\theta} g_{\theta}(z)\|^2}{|1 - D(g_{\theta}(z))|^2} \right]$$

$$|D(x) - D^*(x)| \leq \epsilon \quad \forall x$$

$$\|\nabla_x (D(x) - D^*(x))\| \leq \epsilon \quad \forall x$$

$$\nabla_{\theta} E_z [\log(1 - D(g_{\theta}(z)))] \leq \frac{\epsilon^2}{(1 - \epsilon)^2} E_z [\|\int_{\theta} g_{\theta}(z)\|^2]$$

$$\stackrel{*}{\sim} \frac{\epsilon^2}{(1 - \epsilon)^2} \rightarrow 0 \quad \text{as } \epsilon \rightarrow 0$$

Thm: as $D \rightarrow D^*$, $\nabla_{\theta} L(D, g_{\theta}) \rightarrow 0$

Not covered: $\Delta \theta = \nabla_{\theta} E_z [-\log D(g_{\theta}(z))]$

$\uparrow \downarrow$; \rightarrow Cauchy dist. (under general assumptions)

3. Random Noise \cup

Replace P_g and P_r with $P_{g+\epsilon}$, $P_{r+\epsilon}$

$$P_{x+\epsilon}(x) = E_{y \sim P_x} [P_{\epsilon}(x-y)] = \int P_{\epsilon}(x-y) dP_x(y)$$

$$\text{Recall } D^*(x) = \frac{P_r(x)}{P_g(x) + P_r(x)} \rightarrow D^*(x) = \frac{P_{r+\epsilon}(x)}{P_{g+\epsilon}(x) + P_{r+\epsilon}(x)}$$

Thm: $E_{z \sim p(z)} [\nabla_{\theta} \log(1 - D^*(g_{\theta}(z)))] \quad \epsilon, \epsilon' \sim N(0, \sigma^2 I)$

$$= E_z \left[\int_{\epsilon'} a(z) \int P_{\epsilon'}(\tilde{g}_{\theta}(z) - y) \nabla_{\theta} \|\tilde{g}_{\theta}(z) - y\|^2 dP_r(y) - b(z) \int P_{\epsilon'}(\tilde{g}_{\theta}(z) - y) \nabla_{\theta} \|\tilde{g}_{\theta}(z) - y\|^2 dP_g(y) \right] = 2 \nabla_{\theta} \text{SSD}(P_{r+\epsilon}, P_{g+\epsilon})$$

$$a(z) = \frac{1}{2\sigma^2} \frac{1}{P_{r+\epsilon}(g_{\theta}(z)) + P_{g+\epsilon}(g_{\theta}(z))}, \quad b(z) = a(z) \frac{P_{r+\epsilon}(g_{\theta}(z))}{P_{g+\epsilon}(g_{\theta}(z))}$$

