

WGAN Loss

$$L = \sum_{x \in \text{real}} h(x) - \sum_{z \in \text{latent}} h(G(z))$$

$$\min_G \max_{h \in \text{Lip}_1} L$$

Enforcing Lipschitz: many approaches

• original paper: "weight clipping"

after each weight update

$$w \rightarrow \text{clip}(w, -c, c) \begin{cases} \text{if } w > c \\ w = c \\ \text{if } w < -c \\ w = -c \end{cases}$$

This enforces K -Lipschitz crudely

roughly: NN output $h = g_1 \circ w_1 \circ g_2 \circ w_2 \circ \dots \circ g_L \circ w_L(x)$

$$\text{So } |\nabla_x h| \sim |\partial_{s_1} v_1 g_2' v_2 \dots g_L' v_L|$$

So, if weights bounded, $|\nabla_x h|$ bounded
& activation gradients g_i' bounded

• Better way: gradient penalty $\lambda \sum_{x \sim p} \|\nabla_x h\|^2$

add regularizer to loss $\mathcal{L} = \mathcal{L} + \lambda \sum_{x \sim p} (\|\nabla_x h\| - 1)^2$

might not be state of the art...

p : dist'n of pts $\hat{x} = tx + (1-t)y$
 $x, y \sim \text{data, gen}$
 $t \sim U(0, 1)$

Some theory behind this...

not clear why seems to require $\|\nabla_x h\| = 1$
instead of ≤ 1 ... but empirically works

- WGAN - MNIST example

- ATLAS Fast colorgan example

- Micheli's slides

Next generative model framework: variational autoencoders

- example of a latent variable model

data $x \rightarrow$ latent variable z

encodes "meaning" of x (aka embedding)
ideally in a much reduced, simpler space

To sample:

$$z \sim p(z)$$

"prior"
some fixed, simple dist'n

$$x \sim p_\theta(x|z)$$

This gives $(x, z) \sim p(x, z)$

Throw away $z \rightarrow x \sim p(x)$.

Objective: max likelihood

want: $p(x) = p_{\text{data}}(x)$.

$$\max_{\theta} \sum_{x \sim \text{data}} \log p_\theta(x) = \max_{\theta} \sum_{x \sim \text{data}} \log \int p_\theta(x|z) p(z) dz$$

↓ direct eval via MCMC?

$$\max_{\theta} \sum_{x \sim \text{data}} \log \frac{1}{N} \sum_{z \sim p(z)} p_\theta(x|z)$$

- Computationally expensive
 - Noisy gradients
 - Curse of dimensionality (MC needs large N in higher dim)
- direct MC generally fails!

Instead: variational approach + tractable lower bound on likelihood

introduce $q_{\beta}(z|x)$ ← proposal dist'n for more efficient sample than prior

$$\log p(x) = \log \int dz p_0(x|z) p(z) q_{\beta}(z|x)$$

$$\approx \log \sum_{z \sim q_{\beta}(z|x)} \frac{p_0(x|z) p(z)}{q_{\beta}(z|x)} \quad \left(\text{sample } z \text{ from } q_{\beta}(z|x) \right)$$

→ sample would be most efficient if

" Jensen
ineq "

$$\geq \frac{1}{N} \sum_{z \sim q_{\beta}(z|x)} \log \frac{p_0(x|z) p(z)}{q_{\beta}(z|x)}$$

$$\log \sum f \geq \sum \log f$$

(AKA AM-GM
ineq)

"Evidence Lowerbound"

ELBO

the inside log wants to be $p(x)$

Bayesian evidence

$q_{\beta}(z|x) \sim p_0(x|z) p(z)$
i.e. integrand of (1)
always

So q_{β} wants to be posterior!
more to come...

The ELBO is tractable!

$$\max_{\theta, \beta} \sum_{x \sim p_{\text{data}}(x)} \text{ELBO} = \max_{\theta, \beta} \sum_{\substack{x \sim p_{\text{data}}(x) \\ z \sim q_{\beta}(z|x)}} \log \frac{p_0(x|z) p(z)}{q_{\beta}(z|x)}$$

Just train w/ batches of x and z together