

Friday, April 2, 2021 3:10 PM

Lecture 19

Problems with GANs

- Convergence: min-max unstable

Vanishing gradients: if discriminator is too powerful, can overwhelm generator
conversely, if discriminator is completely fooled, then just guessing

randomly → gradients no info.

generator cannot learn → generator will start to behave badly.

if GAN achieves good performance

→ quickly destabilize

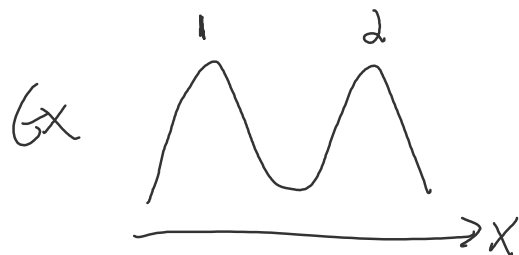
Friday, April 2, 2021 3:28 PM

- Metrics for success

- D & G losses accuracy → not very correlated w/ image quality,
 ↓
 hard to know when to stop training

- Mode collapse

G - D get stuck in a vicious cycle



G can learn to produce very realistic 1's.

D will be fooled for a while.

then D will learn 2's real, guess as 1's.

→ G will switch to making 2's.

Friday, April 2, 2021 3:39 PM

Ways to improve GANs

Many issues w/ GANs related to nature of the loss fn.

$$\text{Vanilla GAN loss } \min_G \max_D L$$

$$\min_G \text{JSD}(p_{data}, p_G)$$

JSD = 1 whenever p_{data} & p_G are disjoint
vanishing gradients!

Vanilla GAN
 loss not very sensitive to diff. btw p_{data} & p_G . (kind of explains mode collapse)

Friday, April 2, 2021 3:44 PM

Wat: measure of similarity of dist's that is unbounded from above,

one promising approach that has been applied to HGP:

Wasserstein GANs (WGANs)

loss based "Earth Movers' Distance" or "Wasserstein Distance" (optimal transport theory)

$$W(P, Q) = \min_{\gamma \in \Pi(P, Q)} \int dx dy \gamma(x, y) \|x - y\| = \min_{\gamma \in \Pi(P, Q)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

$x \in X$
 $y \in X$

γ : joint probability dist'n on $(x, y) \in X \times X$.



$\begin{cases} W=0 \text{ iff } P=Q \\ \text{unbounded from above for } P \neq Q. \end{cases}$

$\Pi(P, Q)$: space of all γ whose marginals are $P(x)$ & $Q(y)$.

$\mathcal{N} \leftrightarrow \mathcal{N} \quad W \rightarrow \infty$
 $\int \gamma = 1$

Friday, April 2, 2021 3:52 PM

W completely intractable — how to use as loss term?

Brilliant idea: use amazing result Kantorovich-Rubinstein duality

$$W(P, Q) = \max_{\|h\|_L \leq 1} \left[E_{x \sim P}[h(x)] - E_{y \sim Q}[h(y)] \right]$$

↑
space of 1-Lipchitz fns

$$|h(x_1) - h(x_2)| \leq \|x_1 - x_2\|$$

K -Lipchitz fn

$$|h(x_1) - h(x_2)| \leq \|x_1 - x_2\| \cdot K$$

$\forall x_1, x_2$ and for some const K .

continuity condition \leftrightarrow equiv. to
saying $|h'(x)| \leq K \quad \forall x$.

parametrise h by a NN!
"critic" network

$$\max_{\|h\|_L \leq 1} (\dots)$$

\hookrightarrow training the critic!

Friday, April 2, 2021 4:06 PM

partial proof of K-R formula

Idea: use Lagrange multipliers

$$L(\gamma, f, g) = \int \|x-y\| \gamma(x,y) dx dy + \int (\rho(x) - \int \gamma(x,y) dy) f(x) dx$$

$$+ \int (\rho(y) - \int \gamma(x,y) dx) g(y) dy.$$

$$= \int_{x \sim p} f(x) + \int_{y \sim q} g(y) + \int dx dy \gamma(x,y) (\|x-y\| - f(x) - g(y))$$

$$W = \min_{\gamma} \max_{f, g} L(\gamma, f, g)$$

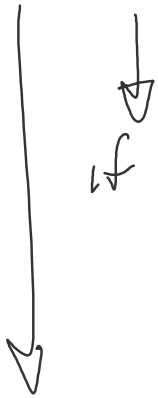
why max?
 f, g

if constraints satisfied \rightarrow Lag. mult. terms give 0
 not satisfied \rightarrow " " give $+\infty$ under max f, g

Friday, April 2, 2021 4:13 PM

Reverse min max order

$$W = \max_{f,g} \min_{\gamma} \left(\int_{x \in \mathcal{X}} [f(x)] + \int_{y \in \mathcal{Q}} [g(y)] + \int dx dy \gamma(x,y) [||x-y|| - f(x) - g(y)] \right)$$



if $f(x) + g(y) \leq ||x-y||$ everywhere then $\min_{\gamma} (\dots) = 0$

$f(x) + g(y) > ||x-y||$ somewhere

$\min_{\gamma} (\dots) = -\infty$ the max spoiled.

$$= \max_{\substack{f,g \\ f(x)+g(y) \leq ||x-y||}} \left(\int_{x \in \mathcal{X}} [f(x)] + \int_{y \in \mathcal{Q}} [g(y)] \right)$$

last step (see literature)

→ prove $\max_{f,g} = \max_{f=-g}$

□.

Friday, April 2, 2021 4:20 PM

Back to WGAN

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

objective

$$\min_G \max_h L$$

\uparrow
h being Lipschitz.

(compare w/ vanilla GAN
 $\min_G \max_D (-BCE)$)

How to enforce Lipschitz constraint?

use weight clipping: after gradient update

clip weights to some fixed range

$w \rightarrow \text{clip}(w, -c, c)$ ↖ hyperparameter.

idea: weight clipping

$h(x; w)$

\hookrightarrow lives in a compact space, so $\nabla_x h$ have to have a maximum somewhere.

Friday, April 2, 2021 4:27 PM

weight clipping is very suboptimal - severely restricting the capacity or expressivity of model.

→ poor performance



better way to implement Lipschitz: Gradient Penalty

$$K\text{-Lipschitz} \Leftrightarrow \max |\nabla_x h| = K.$$

So add regularizer to loss

$$L^* = \lambda \left(\sum_x \left(|\nabla_x h|^2 - 1 \right)^2 \right) \quad \text{gradient penalty term}$$

"WGAN-GP"

(check this)

x is sampled from a mix-mash of $x \sim p_{data}(x)$

$\hat{x} = t x + (1-t) y$ $t \in (0,1)$ $y \sim p_{gen}(y)$ t random from uniform