

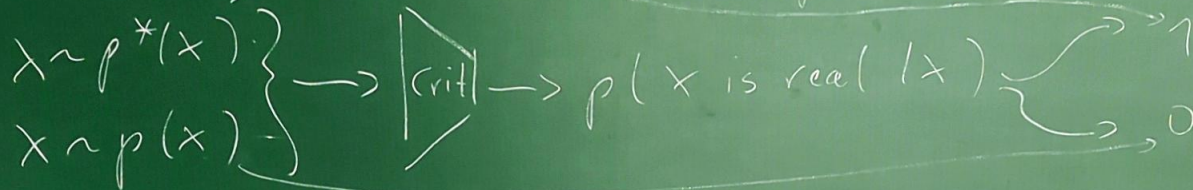
Generative Adversarial Networks (GANs) [Goodfellow et al. 2014]

- dominant generative model 2014 - 2019/20
[now diffusion models, transformers are better]
- idea: learn quality criterion (instead of hand-crafted formula)

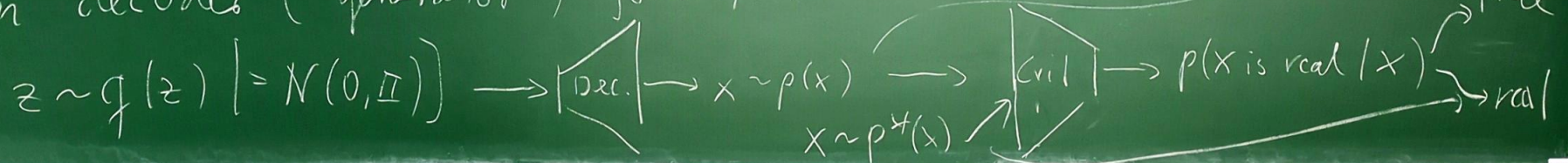
- new neural network: "discriminator" or "critic"

classifier $p(x \text{ is real} | x)$

vs $p(x \text{ is fake} | x) = 1 - p(x \text{ is real} | x)$



- train decoder ("generator") jointly with critic.



- training becomes an "arms race" or "game"
 - critic becomes better at distinguishing reals and fakes
 - decoder becomes better at fooling the critic
- ⇒ training objective

$$\hat{C}, \hat{D} = \underset{D}{\operatorname{argmin}} \underset{C}{\operatorname{argmax}} \left[\mathbb{E}_{x \sim p^*(x)} \left[\log p_C(x \text{ is real} | x) \right] + \mathbb{E}_{z \sim q(z)} \left[\log \left(1 - p_C(\text{real} | D(z)) \right) \right] \right]$$

generated data ↓

$p(x \text{ is fake} | x = D(z))$

- at optimal convergence, we have $p(x) = p^*(x)$

proof: consider expression $a \log y + b \log(1-y)$ with $a, b \neq 0, y \in [0, 1]$
 set derivative $\frac{\partial}{\partial y} \dots = 0$

let $y = p^*(x \text{ is real} | x)$ by the optimal critic, $a = p^*(x), b = p(x)$

$$\text{loss}(D) = \mathbb{E}_{x \sim p^*(x)} \left[\log \frac{p^*(x)}{p^*(x) + p(x)} \right] + \mathbb{E}_{x \sim p(x)} \left[\log \frac{p(x)}{p^*(x) + p(x)} \right]$$

$$= \text{KL} \left[p^*(x) \parallel \frac{p^*(x) + p(x)}{2} \right] + \text{KL} \left[p(x) \parallel \frac{p^*(x) + p(x)}{2} \right] - \log 4$$

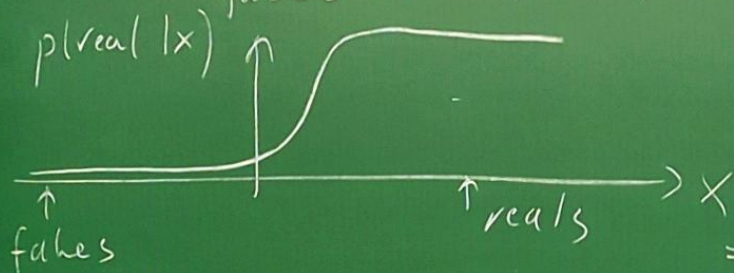
$$\geq -\log 4$$

minimum is achieved if both KL = 0

$$p^*(x) = \frac{p^*(x) + p(x)}{2}$$

$$\text{and } p(x) = \frac{p^*(x) + p(x)}{2} \Leftrightarrow \boxed{p(x) = p^*(x)}$$

- if critic is perfect, $\text{Loss}(D)$ is convex, so gradient descent finds global optimum
- two crucial differences in practice
 - training with perfect critic does not work. for a bad decoder, recognizing fakes is very easy \Rightarrow loss gradient $\approx 0 \Rightarrow$ no training signal



\Rightarrow train D and C jointly, i.e. start with D and C both randomly initialized

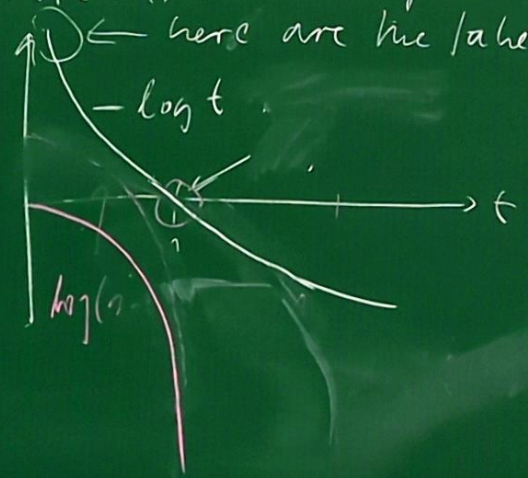
\Rightarrow alternating optimization: 1 step ∇D
 • 1..4 steps ∇C

- use non-saturating loss: replace $\log(1-p(\text{real}(x)))$ with $-\log p(\text{real}(x))$ in decoder training step

$$\hat{C} = \operatorname{argmax}_C \mathbb{E}_{x \sim p^*(x)} [\log p_C(\text{real}(x))] + \mathbb{E}_{z \sim q(z)} [\log(1 - p_C(\text{real}(D(z))))]$$

$$\hat{D} = \operatorname{argmin}_D \mathbb{E}_{z \sim q(z)} [-\log p_C(\text{real}(D(z)))] \quad \text{alternate between the two losses}$$

more informative gradients for D, but the global optimum $p(x) = p^*(x)$ is preserved



why? unclear, hint out!

• GANs behind state-of-the-art up to 2019/20

• variants: - Wasserstein GAN: uses alternative loss for critic but is not really better

- cycle GAN. replace $q(z)$ with true distribution over another variety of real data

$p^*(x)$: distr. of daylight photos, $\tilde{p}^*(\tilde{x})$: distr. night photos
distr. of satellite images maps

- two cases: (1) paired training set: same x from both variants
(2) unpaired training: no overlap between instances

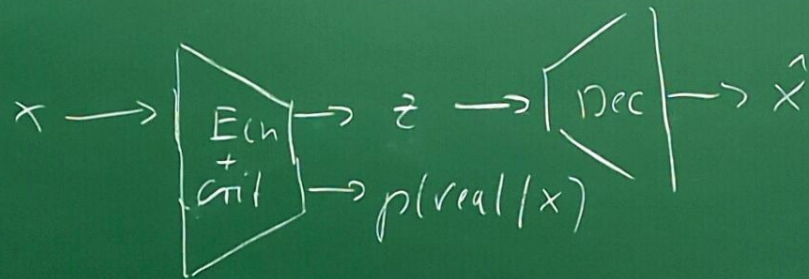
Supervised learning

minimize $\|\tilde{x}_i - E(x_i)\|^2$

unsupervised learning

renew: new "cycle losses": $\|x - D(E(x))\|^2$ and $\|\tilde{x} - E(D(\tilde{x}))\|^2$
should be small both in (1) and (2) optimize quality of instances

- invertible GAN: add an encoder to original GAN



recreating codes $z(x)$ and distinguishing reals from fakes can use same image features

\Rightarrow single network for encodes & critic

\Rightarrow many additional losses as in cycle GAN (esp. cycle losses)

cycles $z \rightarrow \hat{x} \rightarrow \hat{z}$, $x \rightarrow z \rightarrow \hat{x}$ [$x \rightarrow z \rightarrow \hat{x} \rightarrow \hat{z} \rightarrow \hat{\hat{x}}$ etc]
 push forward $E_{\#} p^*(x) \approx q(z)$ self-consistency

critic

$$p_c(\text{real} | x \sim p(x)) \approx p_c(\text{fake} | x \sim p(x))$$

but: hyperparameter optimization is hard

(fakes look like reals)

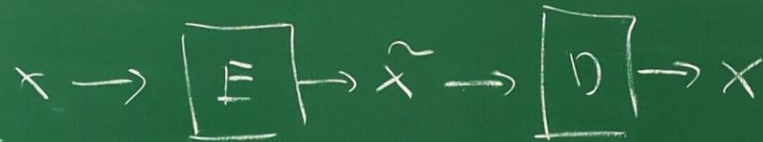
without critic, checking with $D(\tilde{x}) = \tilde{x}$ and $E(x) = x$ would optimize the cycle losses

\Rightarrow critic ensures that $p_E(\tilde{x}) = E_{\#} p^*(x) \approx \tilde{p}^*(\tilde{x})$ } optimize quality of distributions

$$p_D(x) = D_{\#} \tilde{p}^*(\tilde{x}) \approx p^*(x)$$

[if paired TS: two additional losses $\| \tilde{x}_i - E(x_i) \|^2$

difference to autoencoder: no bottleneck $\| x_i - D(\tilde{x}_i) \|^2$)



2 additional losses for training