

PixelRNN/CNN

↳ Fully visible belief network

↳ Decompose likelihood of an image x into product of 1-d distributions using chain rule:

$$p(x) = \prod_{i=1}^n p(z_i | x_1, \dots, x_{i-1})$$

\uparrow

likelihood of image x

probability of i th pixel value given all previous pixels

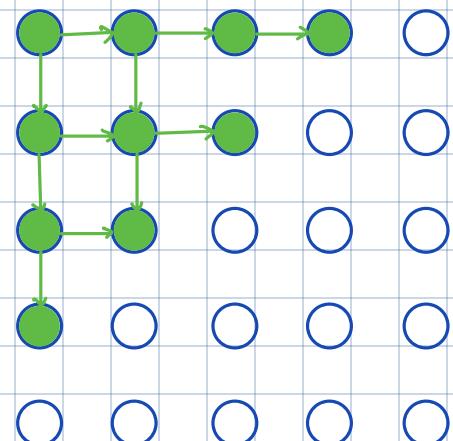
↳ Maximize likelihood of training data

PixelRNN:

↳ Generate image pixels starting from corner

↳ Dependency on previous pixels modeled using RNN (LSTM)

↳ Slow because it is sequential



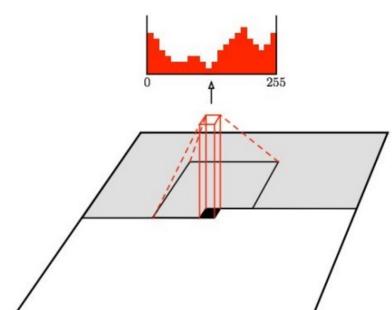
PixelCNN:

↳ Still generate image pixels starting from the corner

↳ Dependency on previous pixels now modeled using a CNN over context region

↳ Training: Maximize likelihood of training images

Softmax loss at each pixel



Variational Autoencoders (VAE)

↳ Define intractable density function with latent z :

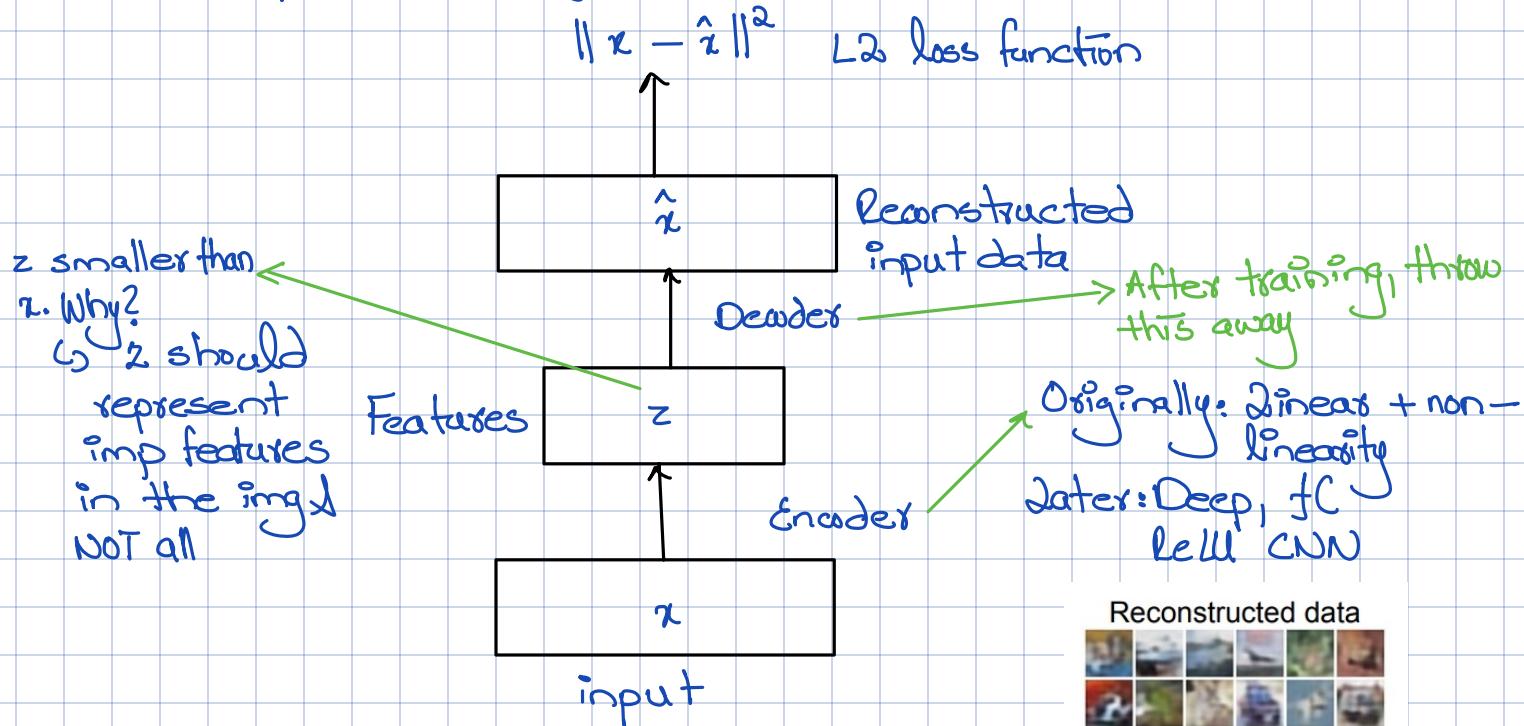
$$P_\theta(x) = \int P_\theta(z) P_\theta(x|z) dz$$

↳ Cannot optimize directly, derive & optimize lower bound on likelihood instead

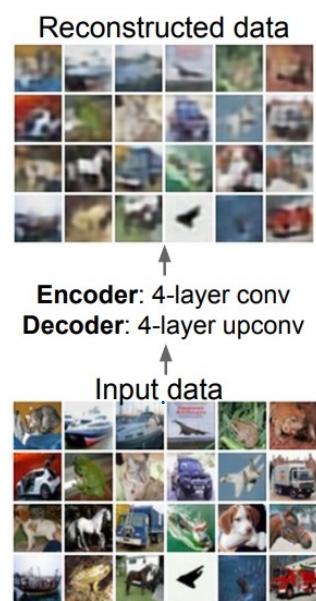
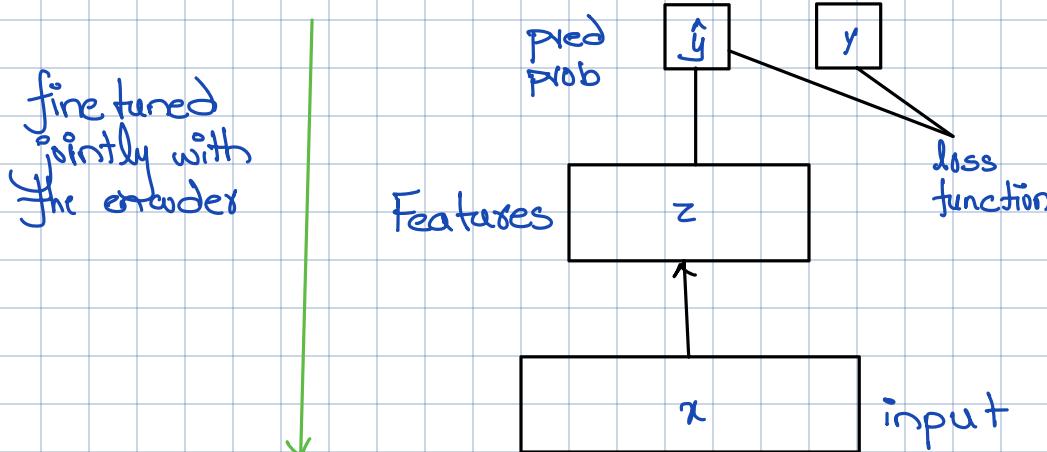
Background:

↳ VAE derived from AE

↳ AE is unsupervised approach for learning a lower dimensional feature representation from unlabeled data.

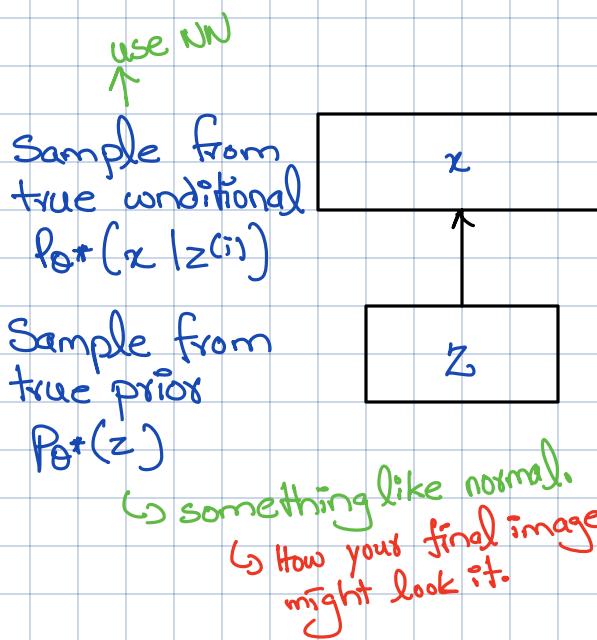


↳ We can use this initialize a supervised model



VAE:

- ↳ Adds probabilistic spin on autoencoders - will let us sample from model to generate data!
- ↳ Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from underlying observed (latent) representation z .



Intuition:

z contains different types of attributes. for example, how much of a smile on the face?

How to train this network?

- ↳ Learn model params to maximize likelihood of training data

$$P_\theta(x) = \int P_\theta(z) P_\theta(x|z) dz$$

computing $P(x|z)$ for every z is intractable

Q. What is the problem?

↳ This is NOT tractable

↳ Complex & computationally expensive

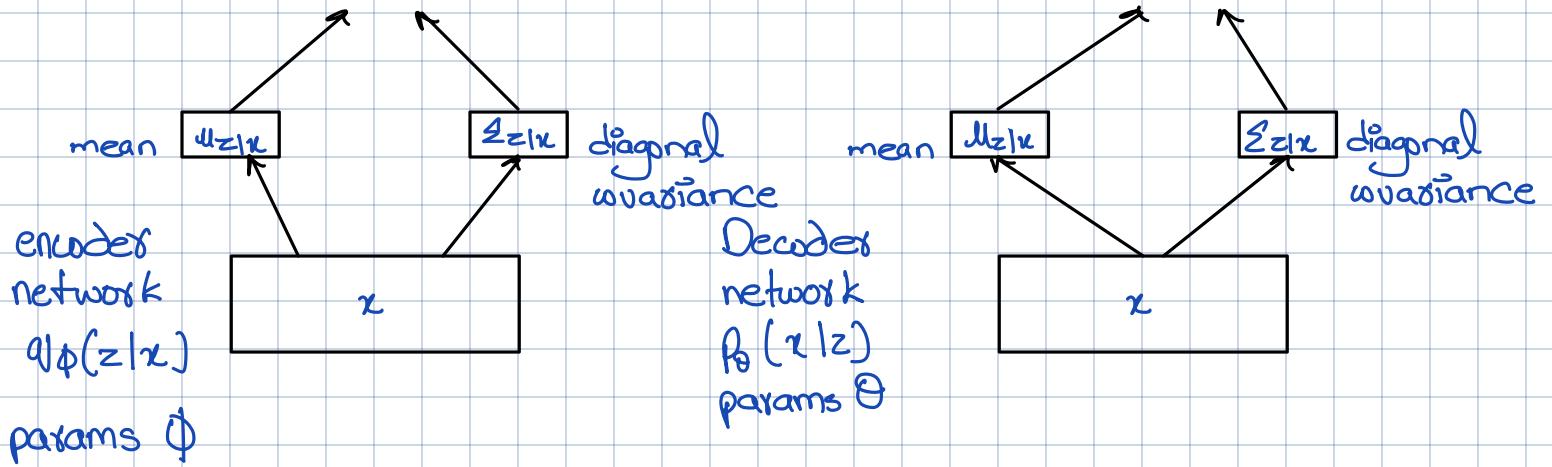
- ↳ Posterior density also intractable: $p(z|x) = \frac{\int P_\theta(x|z) P_\theta(z) dz}{P_\theta(x)}$

Solution:

↳ Use a NN (of course!) called encoder network $q_{\phi}(z|x)$ that approxs $p(z|x)$

↳ Will allow to derive a lower bound on the data likelihood that is tractable, which we can optimize

sample z from $z|x \sim N(\mu_{z|x}, \Sigma_{z|x})$ sample $z|z$ from $z|z \sim N(\mu_{z|z}, \Sigma_{z|z})$



↪ Encoder & decoder network also called recognition / inference & generation networks

Now,

$$\begin{aligned}
 \log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\
 &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\
 &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))
 \end{aligned}$$

Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :(But we know KL divergence always ≥ 0 .

Reconstruct the input data

$$\begin{aligned}
 \log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\
 &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\
 &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{> 0}
 \end{aligned}$$

$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

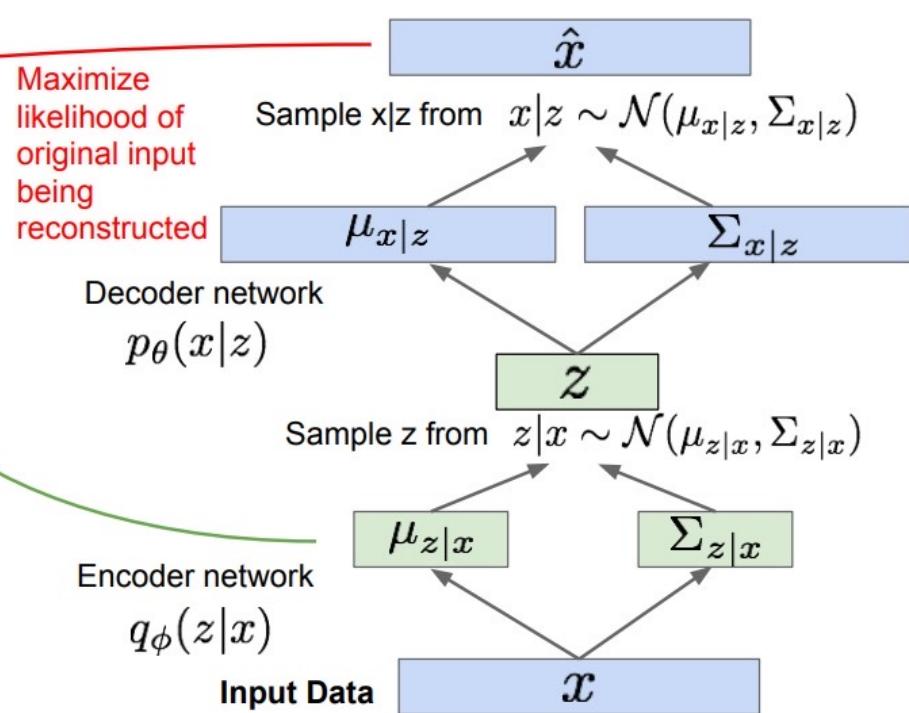
Training: Maximize lower bound

Training data:

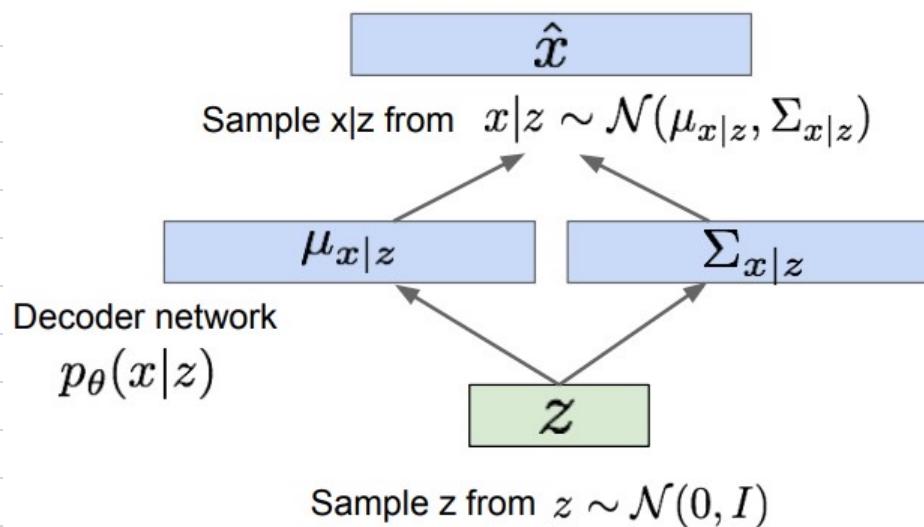
Putting it all together: maximizing the likelihood lower bound

$$\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$\mathcal{L}(x^{(i)}, \theta, \phi)$



Generating data:

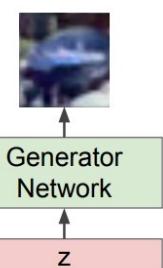


Generative Adversarial networks:

- Pass input of random noise through GAN
- Get an output corresponding to sample from training distribution

Output: Sample from training distribution

Input: Random noise

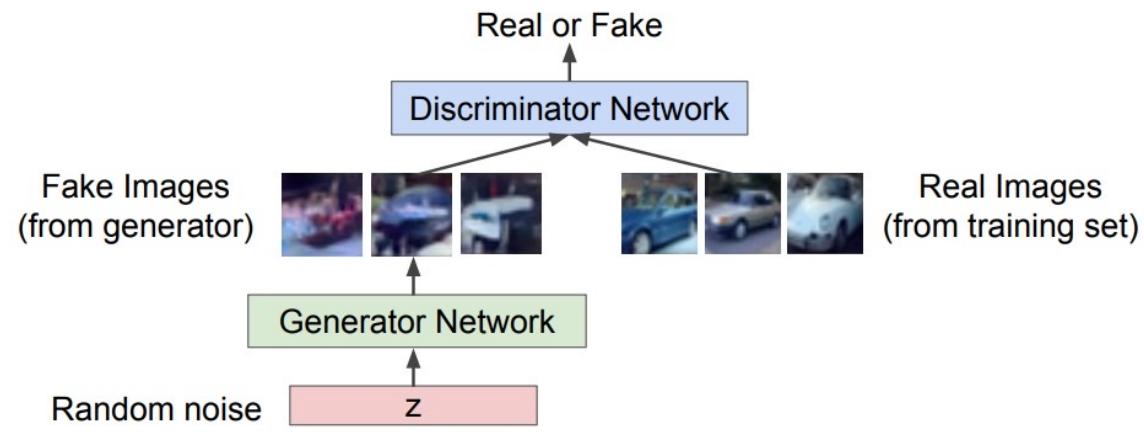


Training GANs:

↳ A two player game

Generators network: Try to fool the discriminator by generating real looking images → Player 1

Discriminator network: Try to distinguish b/w real & fake image → Player 2



↳ Train jointly in a minimax game.

↳ Minimax objective function:

$$\min_{\Theta_g} \max_{\Theta_d} \left[E_{x \sim p_{\text{data}}} \log D_{\Theta_d}(x) + E_{z \sim p(z)} \log (1 - D_{\Theta_d}(G_{\Theta_g}(z))) \right]$$

Discriminator output for real data x
Discriminator output for fake data h(z)

↳ Discriminator outputs likelihood (0,1) of real image

↳ Discriminator wants to maximize objective such that $D(x) \approx 1$ & $D(h(z)) \approx 0$

↳ Generator (Θ_g) wants to minimize objective such that $D(G(z)) \approx 1$

Alternate b/w:

1) Gradient ascent on discriminator:

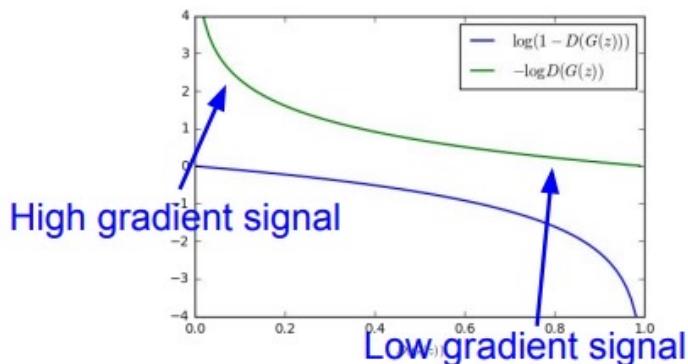
$$\max_{\Theta_d} \left[E_{x \sim p_{\text{data}}} \log D_{\Theta_d}(x) + E_{z \sim p(z)} \log (1 - D_{\Theta_d}(G_{\Theta_g}(z))) \right]$$

2) Gradient descent on generator:

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

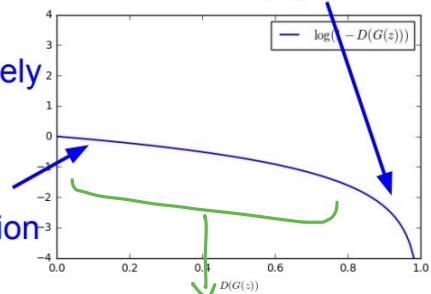
So, we do gradient ascent.

instead of minimizing the probability of the discriminator being correct, we maximize the probability of the discriminator being incorrect.



Gradient signal dominated by region where sample is already good

When sample is likely fake, want to learn from it to improve generator. But gradient in this region is relatively flat!



we want to learn when samples are bad as well

2. Gradient ascent on generator

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} -\log(\text{logit}(z))$$

Algorithm for training GANs:

for number of training iterations do

for k steps do

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(\mathbf{x}^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)}))) \right]$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)})))$$

end for