

(Deep) Generative Models for Unsupervised Learning (Part 4 – GANs)

CS698X: Topics in Probabilistic Modeling and Inference

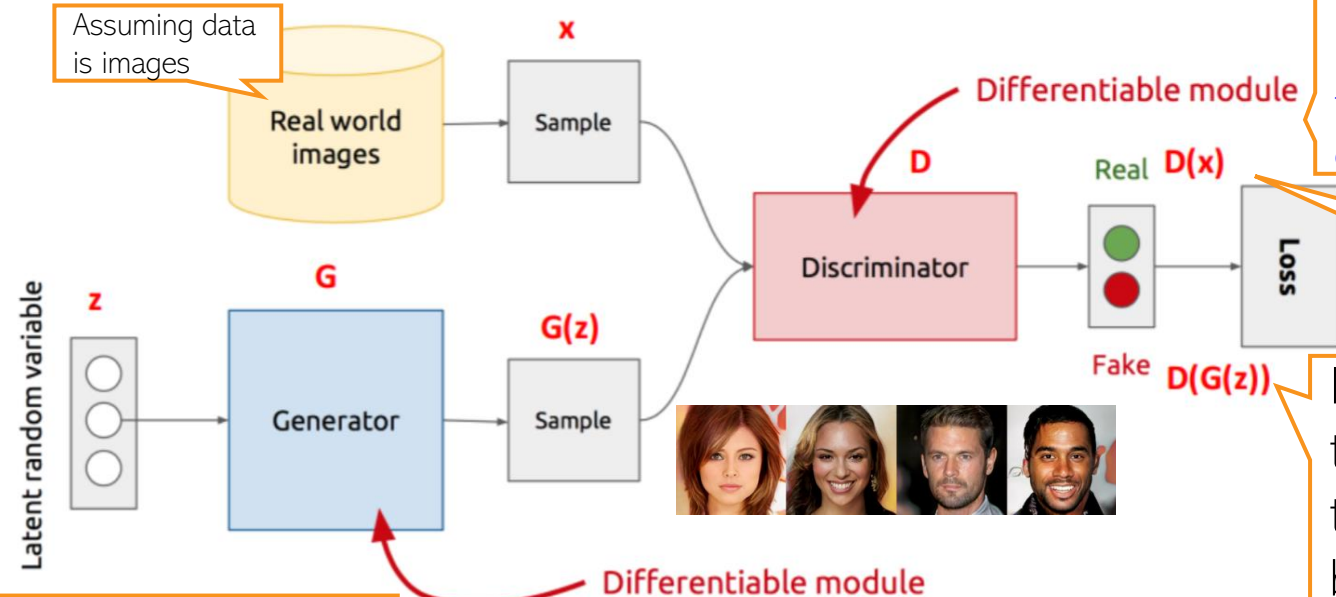
Piyush Rai

Generative Adversarial Network (GAN)

- GAN is an implicit generative latent variable model
- Can generate from it but can't compute $p(\mathbf{x})$ - the model doesn't define it explicitly
- GAN is training using an **adversarial way** (Goodfellow et al, 2013)

Unlike VAE, no explicit parametric likelihood model $p(\mathbf{x}|\mathbf{z})$

Thus can't train using methods that require likelihood (MLE, VI, etc)



The discriminator can be a binary classifier or any method that can compare b/w two distributions (real and fake here)

Discriminator network is trained to make it close to 1

Discriminator network is trained to make it close to 0 and generator network is trained to make it to be close to 1 to fool the discriminator into believing that $G(\mathbf{z})$ is a real sample

Min-max optimization

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Generative Adversarial Network (GAN)

- The GAN training criterion was

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

- With G fixed, the optimal D (exercise)

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

Distribution of real data

Distribution of synthetic data

- Given the optimal D , The optimal generator G is found by minimizing

$$V(D_G^*, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[\log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right]$$

$$= \text{KL} \left[p_{data}(x) \parallel \frac{p_{data}(x) + p_g(x)}{2} \right] + \text{KL} \left[p_g(x) \parallel \frac{p_{data}(x) + p_g(x)}{2} \right] - \log 4$$

Jensen-Shannon
divergence between
 p_{data} and p_g ,
minimized when
 $p_g = p_{data}$

Thus GAN can learn the true data
distribution if the generator and
discriminator have enough modeling power



GAN Optimization

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

- The GAN training procedure can be summarized as

1 Initialize θ_g, θ_d ;

θ_g and θ_d denote the params of the deep neural nets defining the generator and discriminator, respectively

2 **for** *each training iteration* **do**

In practice, for stable training, we run $K > 1$ steps of optimizing w.r.t. D and 1 step of optimizing w.r.t. G

3 **for** K steps **do**

4 Sample minibatch of M noise vectors $\mathbf{z}_m \sim q_z(\mathbf{z})$;

5 Sample minibatch of M examples $\mathbf{x}_m \sim p_D$;

6 Update the discriminator by performing stochastic gradient *ascent* using this gradient:

$$\nabla_{\theta_d} \frac{1}{M} \sum_{m=1}^M [\log D(\mathbf{x}_m) + \log(1 - D(G(\mathbf{z}_m)))] ;$$

7 Sample minibatch of M noise vectors $\mathbf{z}_m \sim q_z(\mathbf{z})$;

8 Update the generator by performing stochastic gradient *descent* using this gradient:

$$\nabla_{\theta_g} \frac{1}{M} \sum_{m=1}^M \log(1 - D(G(\mathbf{z}_m))) ;$$

9 Return θ_g, θ_d

In practice, in this step, instead of minimizing $\log(1 - D(G(\mathbf{z})))$, we **maximize** $\log(D(G(\mathbf{z})))$

Reason: Generator is bad initially so discriminator will always predict correctly initially and $\log(1 - D(G(\mathbf{z})))$ will saturate

Least-Squares GAN (LSGAN)

- Instead of a log-loss objective, another alternative is Least Squares GAN (LSGAN)

- Discriminator minimizes the following loss

Function of
discriminator
network params

$$J^D = \frac{1}{2m} \sum_{i=1}^m \left[(D(\mathbf{x}_i) - 1)^2 \right] + \frac{1}{2m} \sum_{i=1}^m \left[(D(G(\mathbf{z}_i)))^2 \right]$$

m real examples and m fake
examples from the generator

- Generator minimizes the following loss

Function of
generator
network params

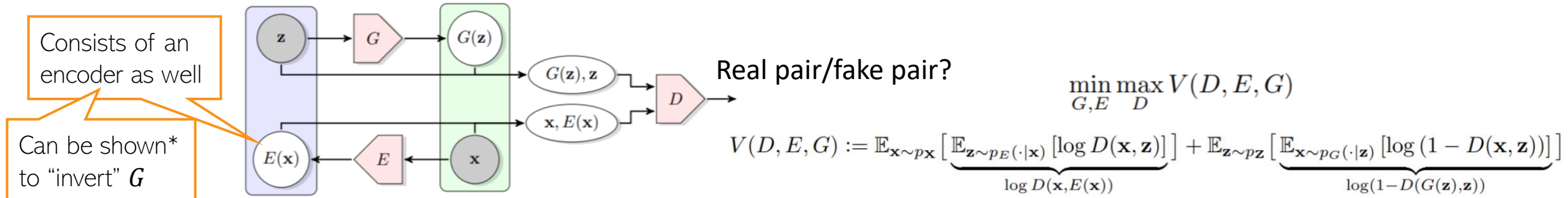
$$J^G = \frac{1}{m} \sum_{i=1}^m \left[(D(G(\mathbf{z}_i)) - 1)^2 \right]$$

- We can perform ALT-OPT for learning the parameters of D and G
- LSGAN address some of the issues, such as saturating gradients

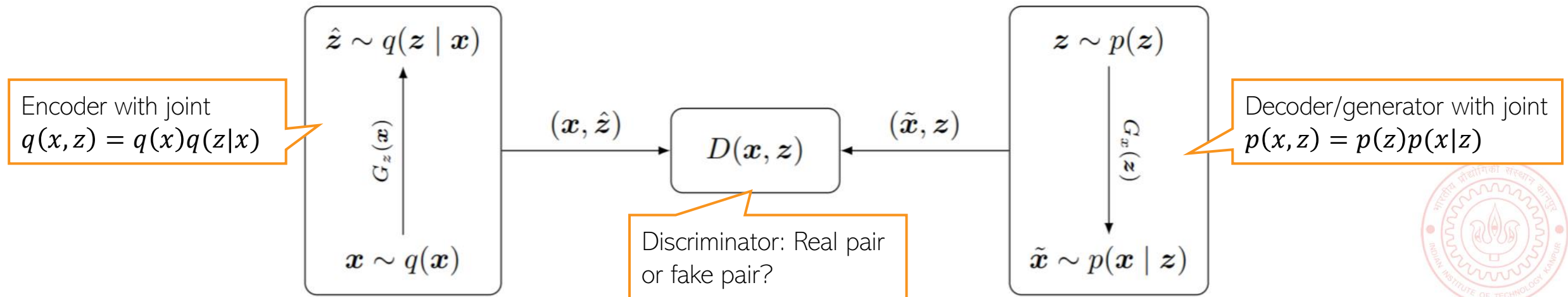


GANs that also learn latent representations

- The standard GAN can only generate data. Can't learn the latent \mathbf{z} from \mathbf{x}
- Bidirectional GAN* (BiGAN) is a GAN variant that allows this



- Adversarially Learned Inference[#] (ALI) is another variant that can learn representations



Evaluating GANs

- Two measures that are commonly used to evaluate GANs
 - Inception score (IS): Evaluates the distribution of generated data
 - Frechet inception distance (FID): Compared the distribution of real data and generated data
- Inception Score will be high (higher is better) if
 - Very few high-probability classes in each sample: Low entropy for its conditional $p(y|x)$
 - We have diverse classes across samples: Marginal $p(y)$ is close to uniform (high entropy)
- FID uses extracted features (using a deep neural net) of real and generated data
 - Usually from the layers closer to the output layer
- These features are used to estimate two Gaussian distributions

Using real data

 $\mathcal{N}(\mu_R, \Sigma_R)$

$\mathcal{N}(\mu_G, \Sigma_G)$

Using generated data
- FID is then defined as
$$\text{FID} = |\mu_G - \mu_R|^2 + \text{trace}(\Sigma_G + \Sigma_R - (\Sigma_G \Sigma_R)^{1/2})$$
 - Lower FID is better

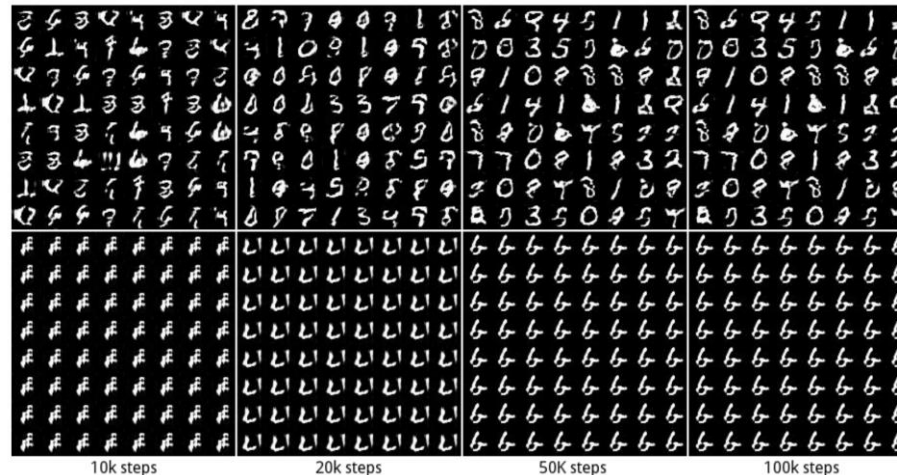
Both IS and FID measure how realistic the generated data is

IS is defined as $\exp(\mathbb{E}_{x \sim p_g} [\text{KL}(p(y|x) || p(y))])$



GAN: Some Issues/Comments

- GAN training can be hard and the basic GAN suffers from several issues
- Instability of training procedure
- Mode Collapse problem: Lack of diversity in generated samples
 - Generator may find some data that can easily fool the discriminator
 - It will stuck at that mode of the data distribution and keep generating data like that



GAN 1: No mode collapse (all 10 modes captured in generation)

GAN 2: Mode collapse (stuck on one of the modes)

- Some work on addressing these issues (e.g., [Wasserstein GAN](#), Least Squares GAN, etc)
- Theoretical properties of GANs yet not fully well-understood (active area of research)



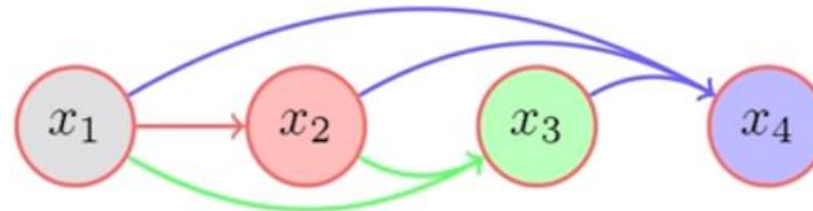
Summary

- Looked at various methods for generative modeling for unsupervised learning
 - Classical methods (FA, PPCA, other latent factor models, topic models, etc)
 - Deep generative models (VAE, GAN)
- Many of these methods can also be extended to model sequential data
- There are also generative models that do not use latent variables
 - Can still be used to generate data and learn the underlying data distribution

Assuming each observation is n-dimensional

$$\prod_{i=1}^n p(x_i | \mathbf{x}_{<k})$$

Can use a neural network to learn (parameters of) each of these distributions



An auto-regressive model

An example: Neural Autoregressive Density Estimator (NADE)

