

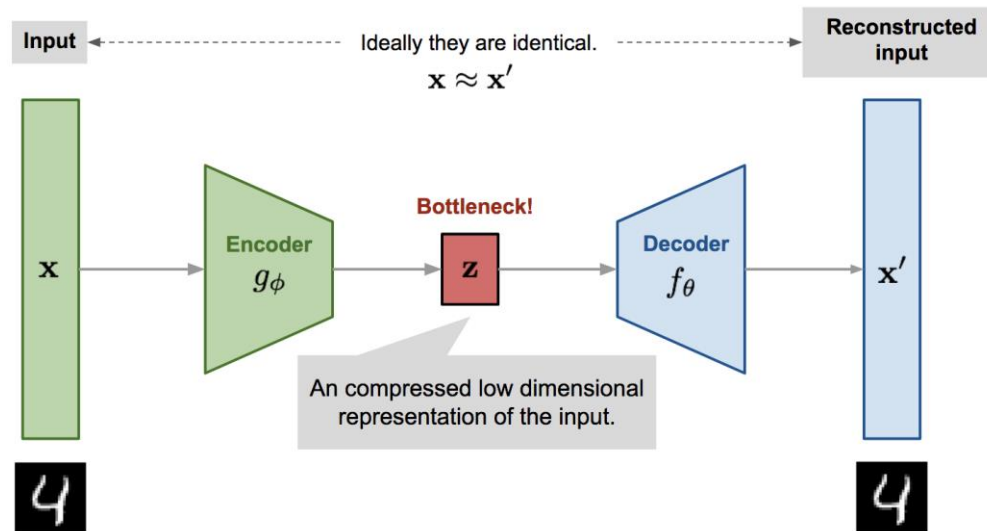
Deep Generative Models (VAE and GAN)

CS772A: Probabilistic Machine Learning

Piyush Rai

A Deep Generative Model: Variational Auto-encoder (VAE)

- VAE* is a probabilistic extension of autoencoders (AE). An AE is shown below



$$L_{AE}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - f_{\theta}(g_{\phi}(\mathbf{x}^{(i)})))^2$$

- The basic difference is that VAE assumes a prior $p(\mathbf{z})$ on the latent code \mathbf{z}
 - This enables it to not just compress the data but also generate synthetic data
 - How: Sample \mathbf{z} from a prior and pass it through the decoder
- Thus VAE can learn good latent representation + generate novel synthetic data
- The name has “Variational” in it since it is learned using VI principles



Variational Autoencoder (VAE)

- VAE has three main components
 - A prior $p_{\theta}(\mathbf{z})$ over latent codes
 - A probabilistic decoder/generator $p_{\theta}(\mathbf{x}|\mathbf{z})$, modeled by a deep neural net
 - A posterior or probabilistic encoder $p_{\theta}(\mathbf{z}|\mathbf{x})$ approx. by an “inference network” $q_{\phi}(\mathbf{z}|\mathbf{x})$

Here θ collectively denotes all the parameters of the prior and likelihood

Using the idea of “Amortized Inference” (next slide)

- VAE is learned by maximizing the ELBO

Here ϕ collectively denotes all the parameters that define the inference network

ELBO for a single data point

$$\begin{aligned}\mathcal{L}(\theta, \phi|\mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathbb{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}))\end{aligned}$$

Maximized to find the optimal θ and ϕ

q_{ϕ} should be such that data \mathbf{x} is reconstructed well from \mathbf{z} (high log-lik)

q_{ϕ} should also be simple (close to the prior)

- The [Reparametrization Trick](#) is commonly used to optimize the ELBO
- Posterior is inferred only over \mathbf{z} , and usually only point estimate on θ



Amortized Inference

- Latent variable models need to infer the posterior $p(\mathbf{z}_n|\mathbf{x}_n)$ for each observation \mathbf{x}_n
- This can be slow if we have lots of observations because
 - We need to iterate over each $p(\mathbf{z}_n|\mathbf{x}_n)$
 - Learning the global parameters needs wait for step 1 to finish for all observations
- One way to address this is via Stochastic VI
- Amortized inference is another appealing alternative (used in VAE and other LVMs too)

$$p(\mathbf{z}_n|\mathbf{x}_n) \approx q(\mathbf{z}_n|\phi_n) = q(\mathbf{z}_n|\text{NN}(\mathbf{x}_n; \mathbf{W}))$$

If q is Gaussian then the NN will output a mean and a variance

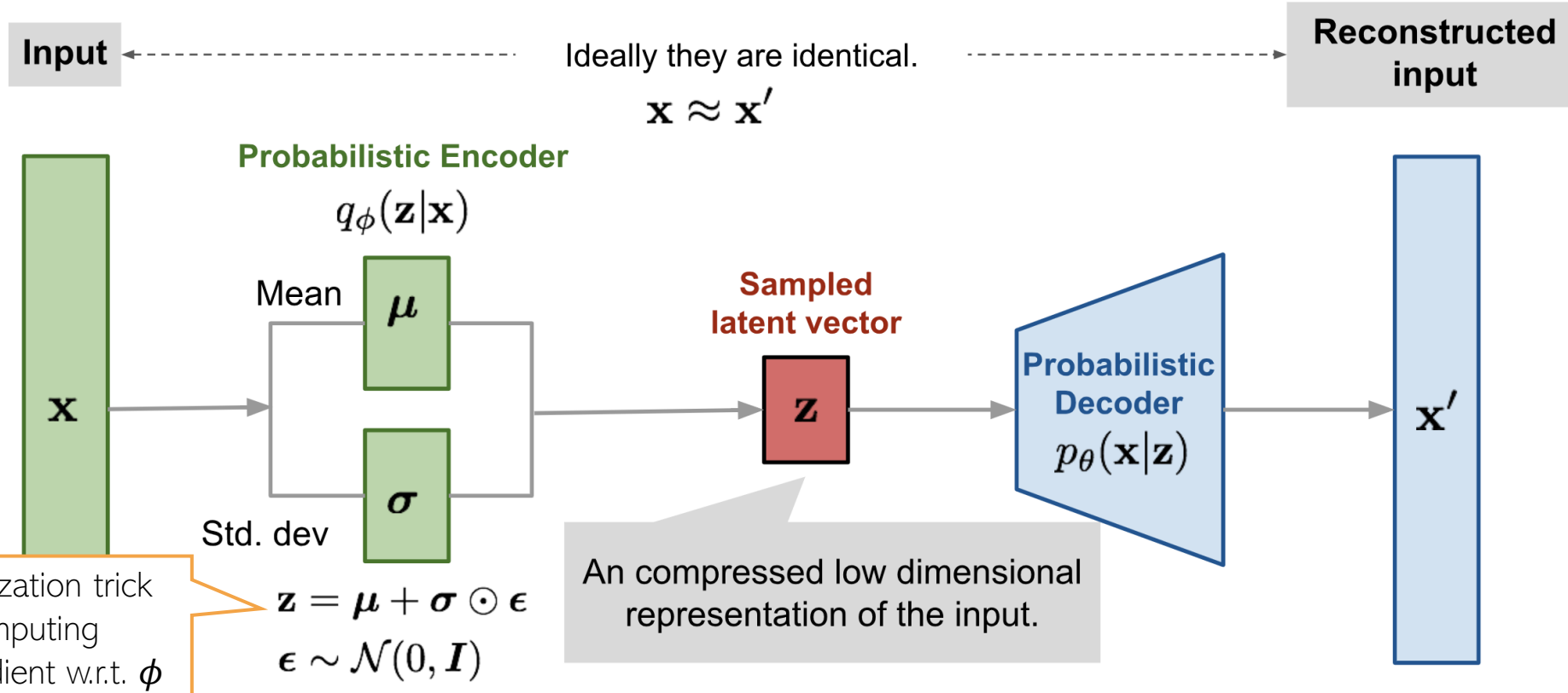
- Thus no need to learn ϕ_n 's (one per data point) but just a single NN with params \mathbf{W}
 - This will be our “encoder network” for learning \mathbf{z}_n
 - Also very efficient to get $p(\mathbf{z}_*|\mathbf{x}_*)$ for a new data point \mathbf{x}_*



Variational Autoencoder: The Complete Pipeline

- Both probabilistic encoder and decoder learned jointly by maximizing the ELBO

$$\begin{aligned}\mathcal{L}(\theta, \phi | \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z} | \mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z}))\end{aligned}$$

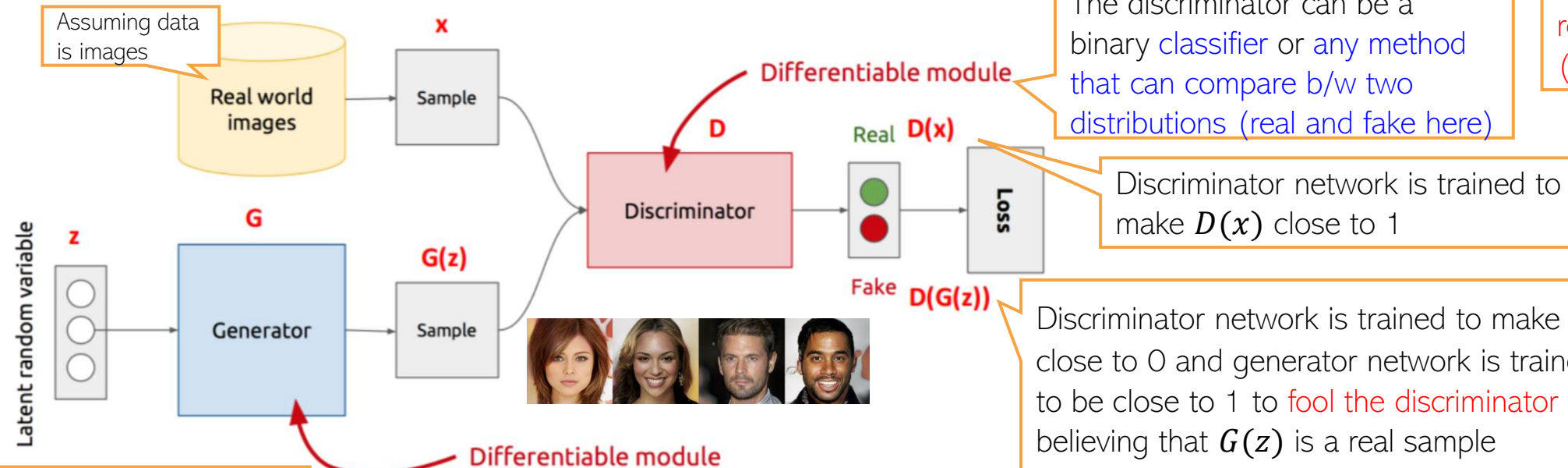


Generative Adversarial Network (GAN)

- 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 trained 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)



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

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

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

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:

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

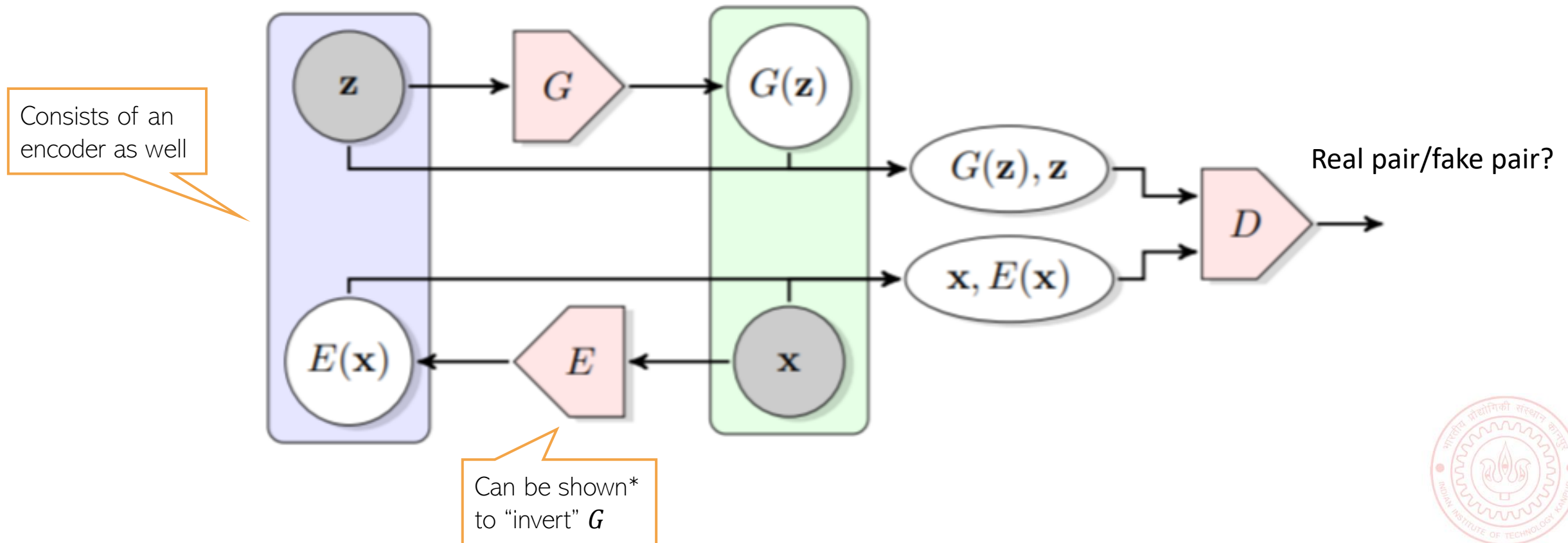
9 Return θ_g, θ_d

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

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



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 defined as $\exp(\mathbb{E}_{x \sim p_g} [\text{KL}(p(y|x) || p(y))])$ will be high if
 - Very few high-probability classes in each sample x : Low entropy for $p(y|x)$
 - 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

High IS and low FID is desirable

Both IS and FID measure how realistic the generated data is

Using real data $\mathcal{N}(\mu_R, \Sigma_R)$

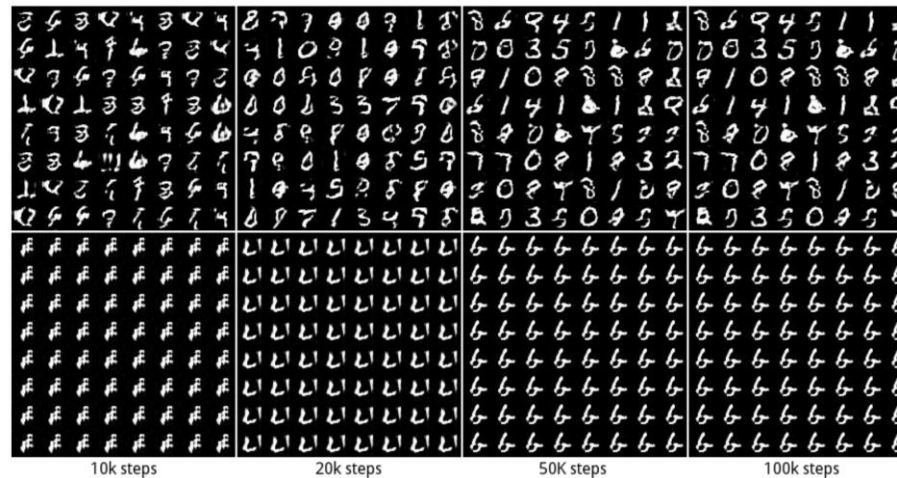
$\mathcal{N}(\mu_G, \Sigma_G)$ Using generated data

- FID is then defined as $\text{FID} = \|\mu_G - \mu_R\|^2 + \text{trace}(\Sigma_G + \Sigma_R - (\Sigma_G \Sigma_R)^{1/2})$
- These measures can also be used for evaluating other deep gen models like VAE



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)

