

# Generative Models

---



GAN

# Generative Adversarial Nets

---

W

# Implicit Generative Model

explicit  $P(x)$

- **Goal:** a sampler  $g(\cdot)$  to generate images
- A simple generator  $g(z; \theta)$ :
  - $z \sim N(0, I)$
  - $x = g(z; \theta)$  deterministic transformation

$\{x_1, \dots, x_N\}$

- Likelihood-free training:
  - Given a dataset from some distribution  $P_{data}$
  - Goal:  $g(z; \theta)$  defines a distribution, we want this distribution  $\approx P_{data}$
  - Training: minimize  $D(g(z; \theta), P_{data})$ 
    - $D$  is some distance metric (not likelihood)
  - Key idea: **Learn a differentiable  $D$**

Metrics:  
① Kullback-Leibler divergence (KL-divergence)  
② Total variation  
③ Wasserstein distance  
- - -

# GAN (Goodfellow et al., '14)

- Parameterize the discriminator  $D(\cdot; \phi)$  with parameter  $\phi$
- **Goal:** learn  $\phi$  such that  $D(x; \phi)$  measures how likely  $x$  is from  $p_{data}$ 
  - $D(x, \phi) = 1$  if  $x \sim p_{data}$
  - $D(x, \phi) = 0$  if  $x \not\sim p_{data}$
  - a.k.a., a binary classifier
- GAN: use a neural network for  $D(\cdot; \phi)$
- **Training:** need both negative and positive samples
  - Positive samples: just the training data  $\{x_1, \dots, x_N\}$
  - Negative samples: use our sampler  $g(\cdot; z)$  (can provide infinite samples).

$$D(\cdot, \cdot) \in [0, 1]$$

- **Overall objectives:**

- Generator:  $\theta^* = \max_{\theta} D(g(z; \theta); \phi)$

- Discriminator uses MLE Training:

$$\phi^* = \max_{\phi} \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim g(\cdot)} [\log(1 - D(\hat{x}; \phi))]$$

# GAN (Goodfellow et al., '14)

- Generator  $G(z; \theta)$  where  $z \sim N(0, I)$ 
  - Generate realistic data
- Discriminator  $D(x; \phi)$ 
  - Classify whether the data is real (from  $p_{data}$ ) or fake (from  $G$ )

- Objective function:

$$L(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim G} [\log(1 - D(\hat{x}; \phi))]$$

- Training procedure:

- Collect dataset  $\{(x, 1) \mid x \sim p_{data}\} \cup \{(\hat{x}, 0) \sim g(z; \theta)\}$

- Train discriminator

$$D : L(\phi) = \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim G} [\log(1 - D(\hat{x}; \phi))]$$

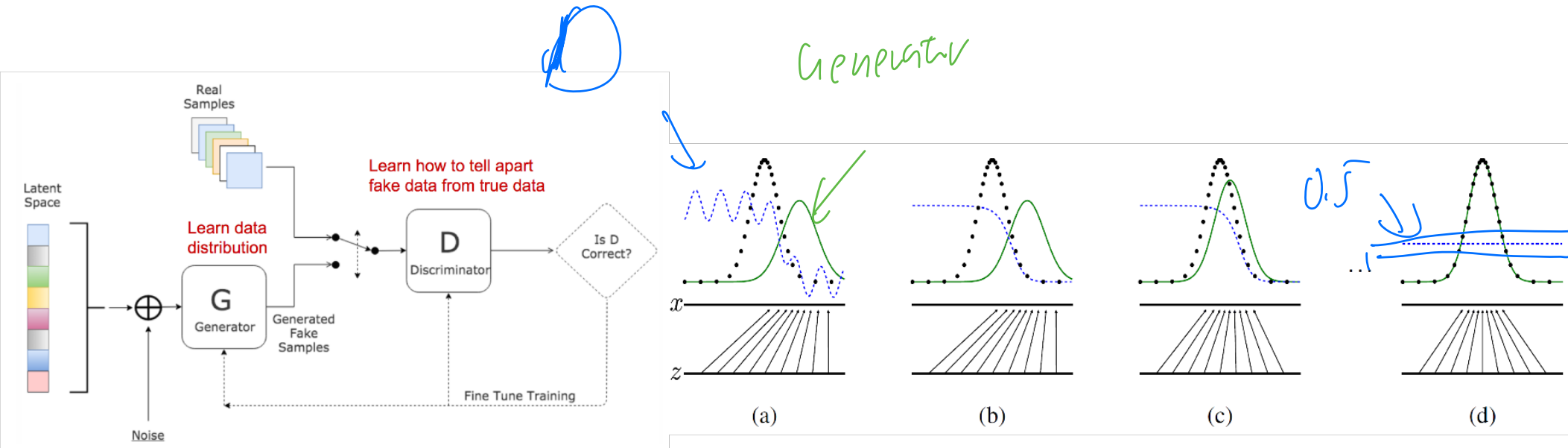
- Train generator  $G : L(\theta) = \mathbb{E}_{z \sim N(0, I)} [\log D(G(z; \theta), \phi)]$

- Repeat

# GAN (Goodfellow et al., '14)

- Objective function:

$$L(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim G} [\log(1 - D(\hat{x}; \phi))]$$



# Math Behind GAN

$$\mathcal{L}(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim p_g(\cdot)} [\log(1 - D(\hat{x}; \phi))]$$

Let  $D^*$ ,  $g^*$  be the solution to  $\mathcal{J}$

Optimal  $D^*$ , for a given  $x$

$$\mathcal{L}_x(D) = p_{\text{data}}(x) \cdot \log D(x) + p_g(x) \cdot \log(1 - D(x))$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial D} = 0 \quad , \quad \text{first-order condition}$$

$$\Rightarrow \frac{p_{\text{data}}(x)}{D^*(x)} - \frac{p_g(x)}{1 - D^*(x)} = 0$$

$$\Rightarrow D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

(if  $p_{\text{data}} = p_g$ ,  $D^*(x) = 0.5$ )

# Math Behind GAN

Consider optimal generator  $g^*$ , given optimal  $D^*$

$$L(\theta, \phi) = \mathbb{E}_{x \sim P_{data}} \left[ \log \frac{P_{data}(x)}{P_{data}(x) + P_g(x)} \right] + \mathbb{E}_{x \sim g} \left[ \log \frac{P_g(x)}{P_{data}(x) + P_g(x)} \right]$$

(will show  $g^* = P_{data}$ )

$$= \mathbb{E}_{x \sim P_{data}} \left[ \log \left( \frac{P_{data}(x)}{P_{data}(x) + P_g(x)} \right) \right] - \log 2$$

$$+ \mathbb{E}_{x \sim g} \left[ \log \left( \frac{P_g(x)}{P_{data}(x) + P_g(x)} \right) \right] - \log 2$$

KL(P||Q)  
=  $\mathbb{E}_P \left[ \log \frac{P}{Q} \right]$

$$= \text{KL} \left( P_{data} \parallel \frac{1}{2} (P_{data} + P_g) \right) + \text{KL} \left( P_g \parallel \frac{1}{2} (P_{data} + P_g) \right) - \log 4$$

2. Jensen-Shannon Divergence (JSD)



# KL-Divergence and JS-Divergence

$$\rightarrow K(p \parallel q) = \mathbb{E}_{x \sim p} \left[ \log \frac{p(x)}{q(x)} \right]$$

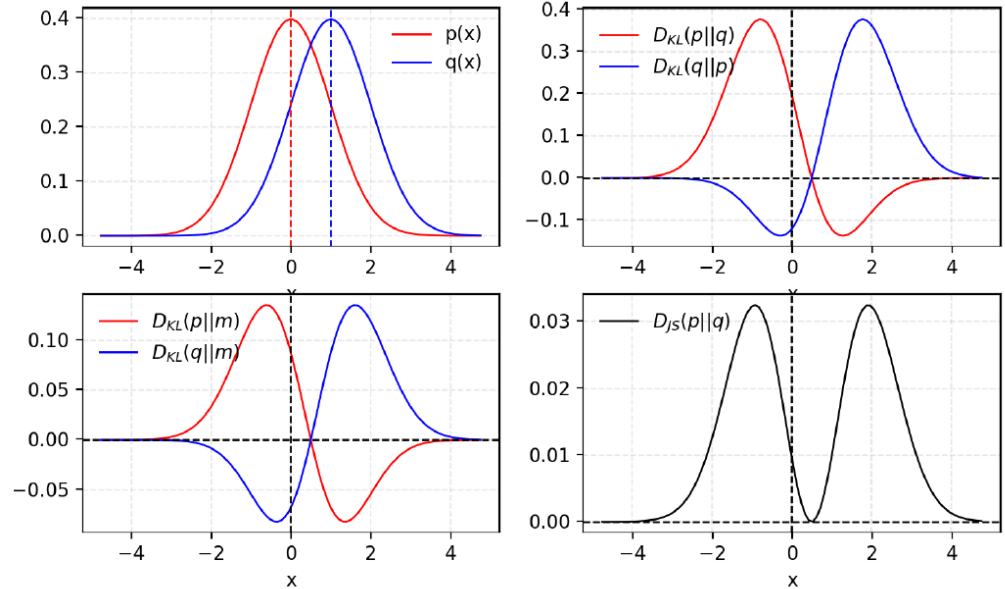
→ not symmetric

$$\rightarrow \text{JSD}(p \parallel q) = \frac{1}{2} \left( K(p \parallel \frac{p+q}{2}) + K(q \parallel \frac{p+q}{2}) \right)$$

→ symmetric

$$\text{JSD} \geq 0$$

$$\text{JSD}(p \parallel q) = 0 \Leftrightarrow p = q$$



# Math Behind GAN

$\Rightarrow$  Given  $D^*$

$$\min_g L(g) = 2 \text{JSD}(P_g \| P_{\text{data}}) - \log 4$$

$\Rightarrow$  global minimizer  $g^*$  :

$$g^* = P_{\text{data}}$$

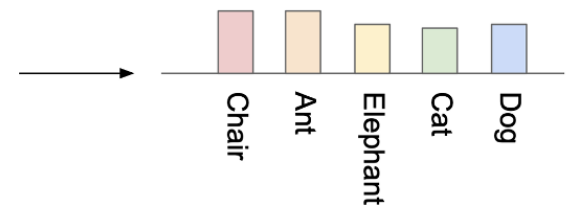
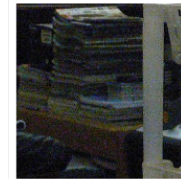
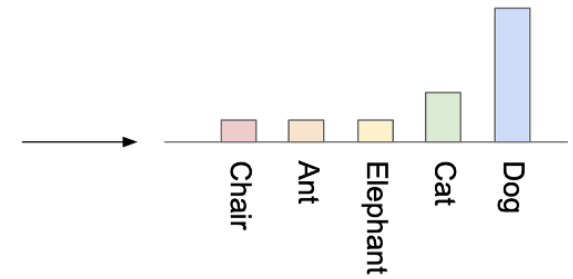
□

# Evaluation of GAN

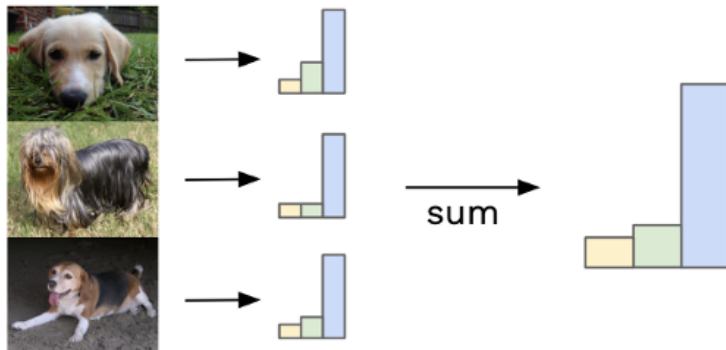
$$f(y|x) \in \mathcal{L}^K$$

- No  $p(x)$  in GAN.
- Idea: use a trained classifier  $f(y|x)$ :
- If  $x \sim p_{data}$ ,  $f(y|x)$  should have low entropy
  - Otherwise,  $f(y|x)$  close to uniform.
- Samples from  $G$  should be diverse:
  - $p_f(y) = \mathbb{E}_{x \sim G}[f(y|x)]$  close to uniform.

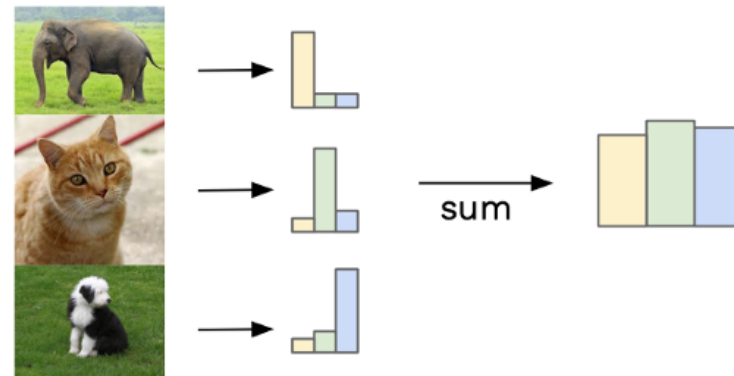
$$K\text{-classes} \in \mathcal{L}^K$$



Similar labels sum to give focussed distribution



Different labels sum to give uniform distribution



# Evaluation of GAN

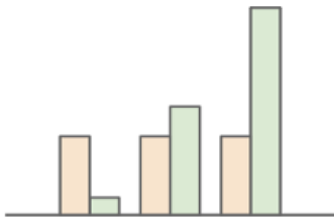
- **Inception Score (IS, Salimans et al. '16)**

- Use Inception V3 trained on ImageNet as  $f(y|x)$

- $IS = \exp \left( \mathbb{E}_{x \sim G} \left[ \underbrace{KL(f(y|x) || p_f(y))}_{\text{one-sample}} \right] \right)$  *marginal*

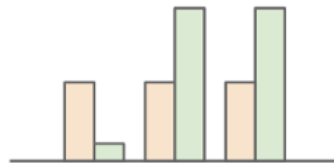
- Higher the better

High KL divergence



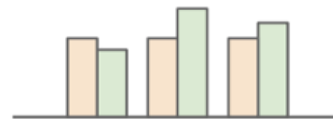
Ideal situation

Medium KL divergence



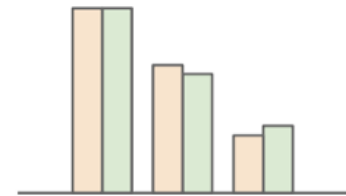
Generated images are not distinctly one label

Low KL divergence



Generated images are not distinctly one label

Low KL divergence



Generator lacks diversity

Label distribution

Marginal distribution

# Comments on GAN

---

- Other evaluation metrics:
  - Fréchet Inception Distance (FID): Wasserstein distance between Gaussians
- Mode collapse:
  - The generator only generate a few type of samples.
  - Or keep oscillating over a few modes.
- Training instability:
  - Discriminator and generator may keep oscillating
  - Example:  $-xy$ , generator  $x$ , discriminatory. NE:  $x = y = 0$  but GD oscillates.
  - No stopping criteria.
  - Use Wasserstein GAN (Arjovsky et al. '17):
$$\min_G \max_{f: \text{Lip}(f) \leq 1} \mathbb{E}_{x \sim p_{data}} [f(x)] - \mathbb{E}_{\hat{x} \sim p_G} [f(\hat{x})]$$
  - And need many other tricks...

# Variational Autoencoder

---



Auto-encoder

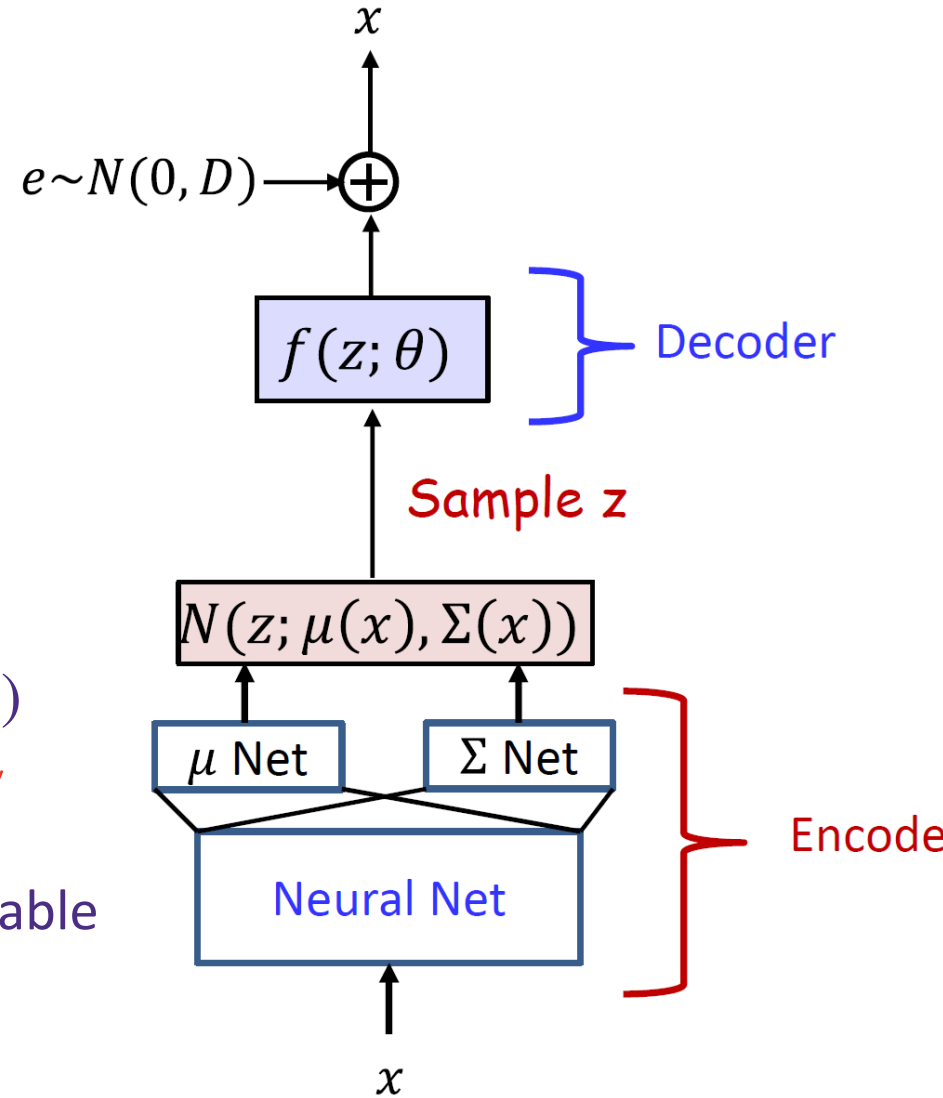
$$D(E(x)) \approx x$$

# Architecture

- Auto-encoder:  $x \rightarrow z \rightarrow x$
- Encoder:  $q(z|x; \phi) : x \rightarrow z$
- Decoder:  $p(x|z; \theta) : z \rightarrow x$

- Isomorphic Gaussian:  $q(z|x; \phi) = N(\mu(x; \phi), \text{diag}(\exp(\sigma(x; \phi))))$
- Gaussian prior:  $p(z) = N(0, I)$
- Gaussian likelihood:  $p(x|z; \theta) \sim N(f(z; \theta), I)$

- Probabilistic model interpretation: latent variable model.



# VAE Training

← induced from variational inference

- Training via optimizing ELBO

- $L(\phi, \theta; x) = \mathbb{E}_{z \sim q(z|x; \phi)} [\log p(z|x; \theta)] - \underbrace{KL(q(z|x; \phi) || p(z))}_{\text{KL penalty}}$
- Likelihood term + KL penalty

- KL penalty for Gaussians has closed form.

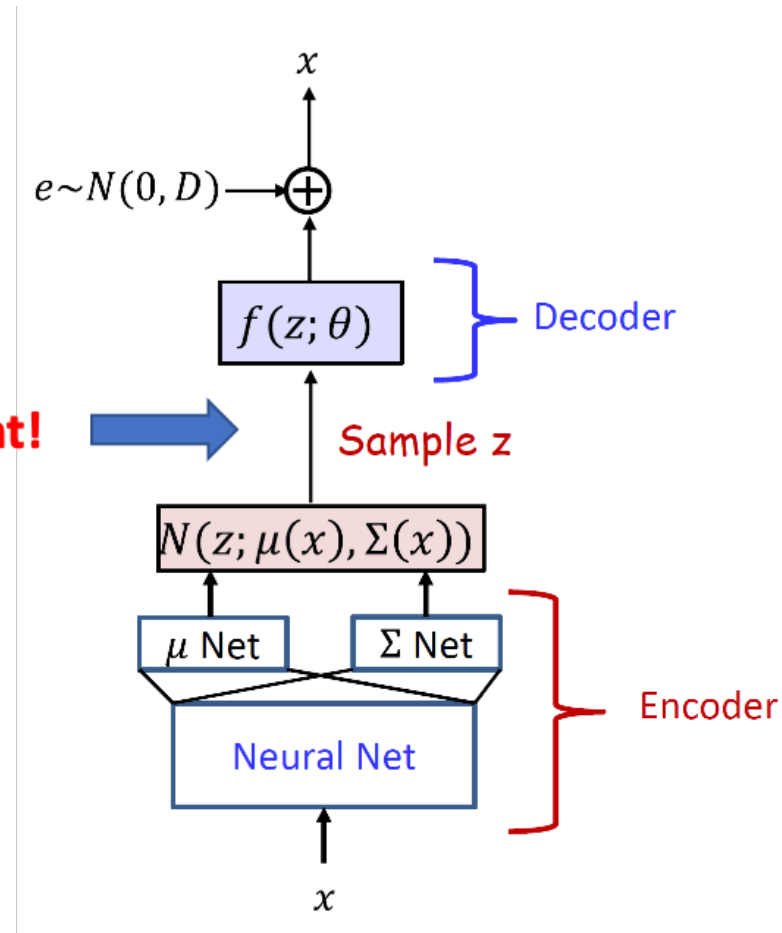
- Likelihood term (reconstruction loss):

- Monte-Carlo estimation  $\{z_1, \dots, z_N\}$
- Draw samples from  $q(z|x; \phi)$
- Compute gradient of  $\theta$ :

- $x \sim N(f(z; \theta); I)$

- $p(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2} \|x - f(z; \theta)\|_2^2)$

$\{x_1, \dots, x_n\}$





# VAE Training

- Likelihood term (reconstruction loss):

- Gradient for  $\phi$ . Loss:  $L(\phi) = \mathbb{E}_{z \sim q(z; \phi)} [\log p(x | z)]$

- Reparameterization trick:

- $z \sim N(\mu, \Sigma) \Leftrightarrow z = \mu + \epsilon, \epsilon \sim N(0, \Sigma)$   $\Leftrightarrow z = \mu + \sigma \cdot \epsilon'$

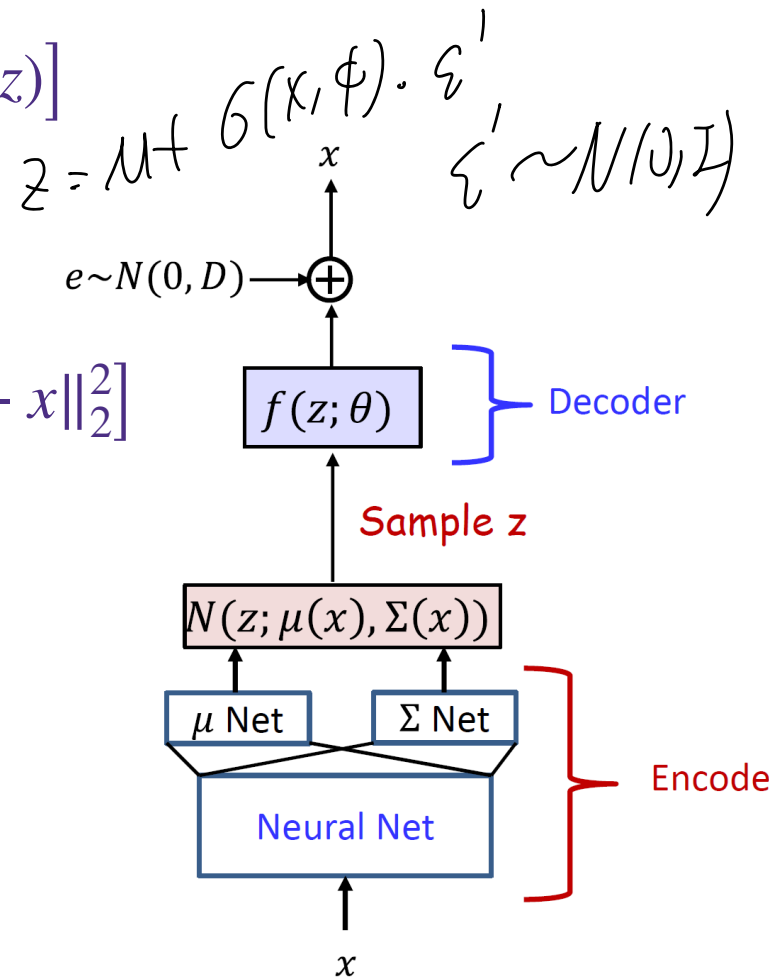
- $L(\phi) \propto \mathbb{E}_{z \sim q(z; \phi)} [\|f(z; \theta) - x\|_2^2]$

- $\propto \mathbb{E}_{\epsilon \sim N(0, I)} [\|f(\mu(x; \phi) + \sigma(x; \phi) \cdot \epsilon; \theta) - x\|_2^2]$

- Monte-Carlo estimate for  $\nabla L(\phi)$

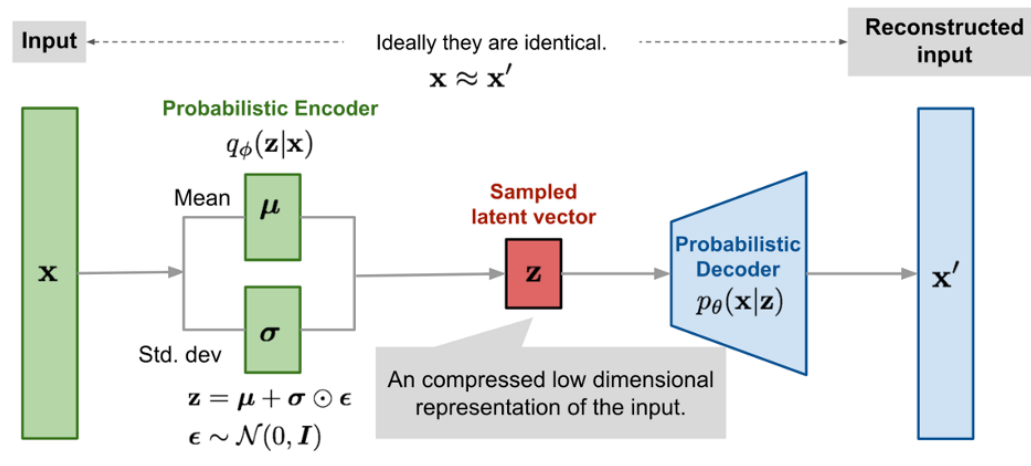
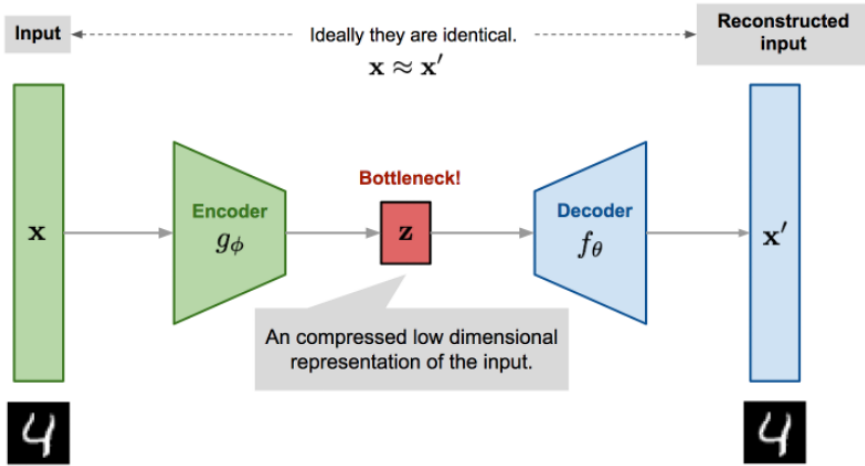
Given  $x$ , Sample  $\epsilon \sim N(0, I)$

- End-to-end training



# VAE vs. AE

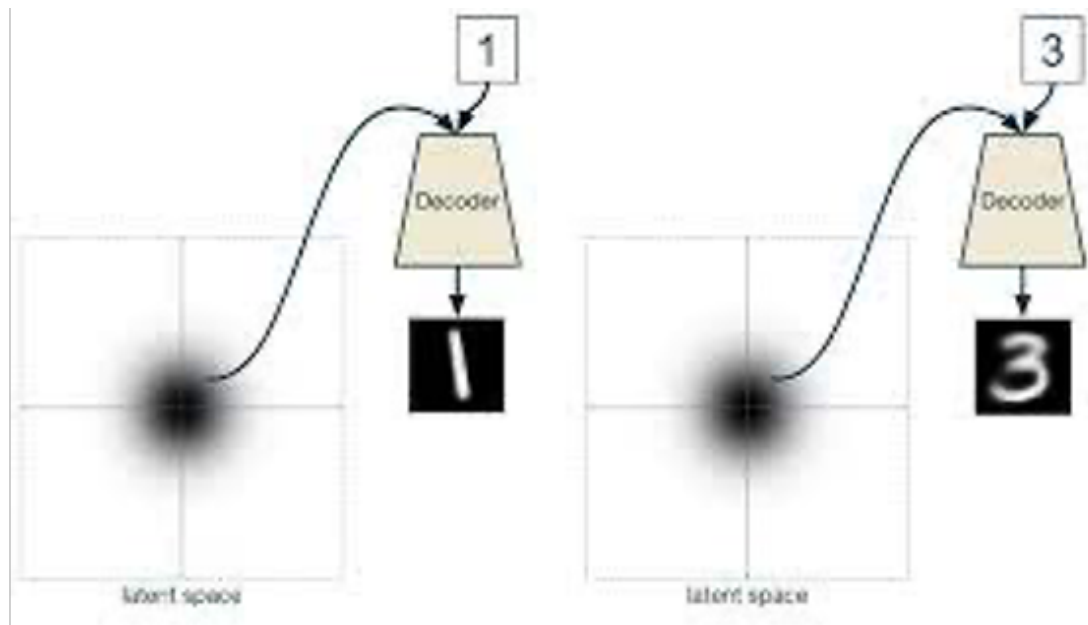
- AE: classical unsupervised representation learning method.
- VAE: a probabilistic model of AE
  - AE + Gaussian noise on  $z$
  - KL penalty:  $L_2$  constraint on the latent vector  $z$



# Conditioned VAE

---

- Semi-supervised learning: some labels are also available



conditioned generation

# Comments on VAE

---

- Pros:
  - Flexible architecture
  - Stable training
- Cons:
  - Inaccurate probability evaluation (approximate inference)

$$p(y)$$