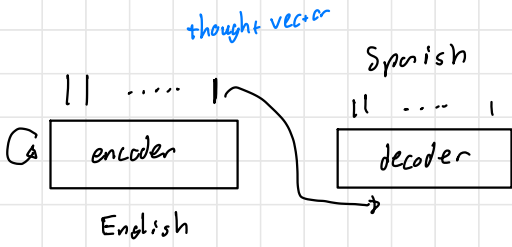


Lecture 19

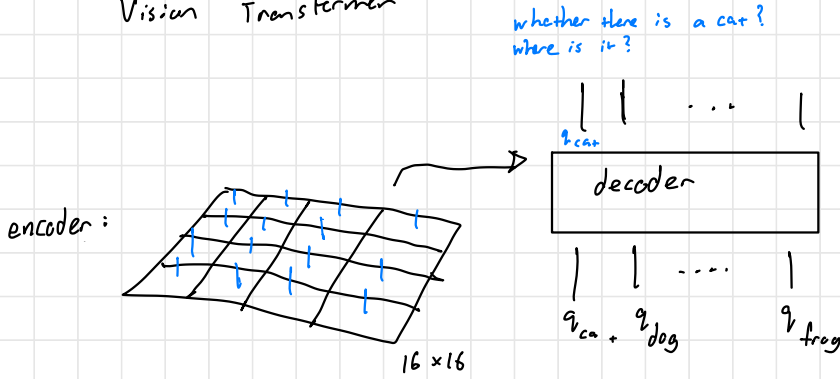


- Neural Network \equiv team of vectors.

Transfer

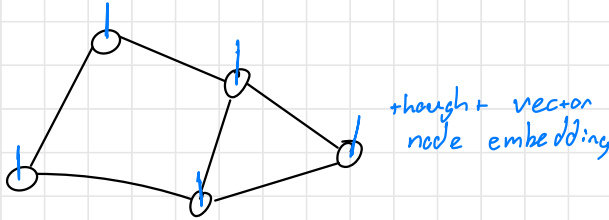


Vision Transformer

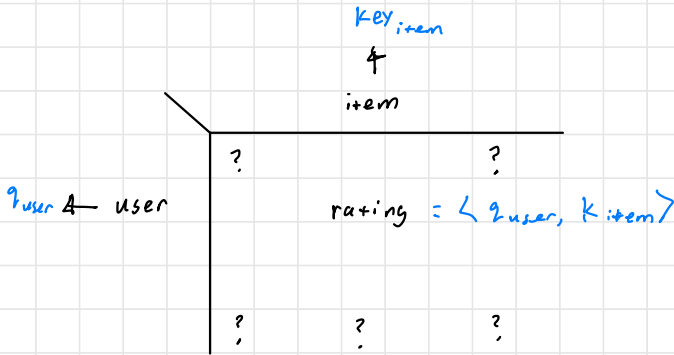


• Graph Neural Network (GNN)

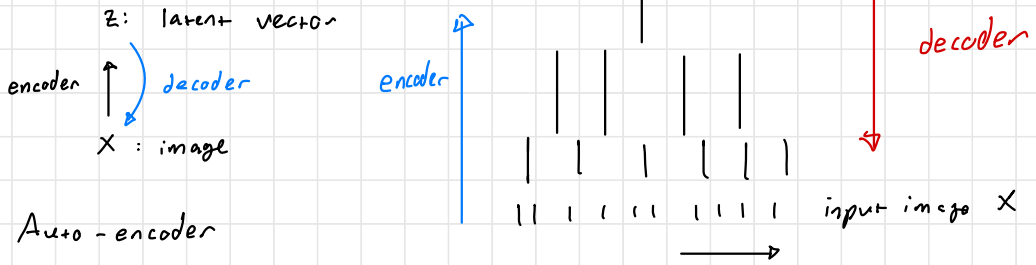
- social network



• Recommender System:



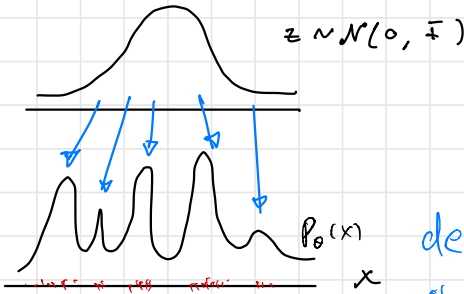
Generative Model: **unsupervised learning**



Generator Model (decoder)

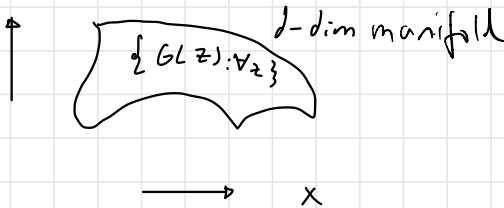
$$z \sim \mathcal{N}(0, I_d)$$

$$X = G(z) + \epsilon$$



density estimation
observed examples sparse in high dimensional space

Manifold Principle:



- GAN (Generative Adversarial Networks)

- Discriminator: $P(x) = P(x \text{ is real}) = P(y=1 | x)$

Real	\vdots i \vdots n	x_i	$y_i = 1$
Fake	\vdots i \vdots n	$\tilde{x}_i = G(\tilde{z}_i)$	$y_i^* = 0$

- log likelihood: $\frac{1}{n} \sum_{i=1}^n \log P(y_i | x_i) + \frac{1}{n} \sum_{i=1}^n \log P(\tilde{y}_i | \tilde{x}_i)$

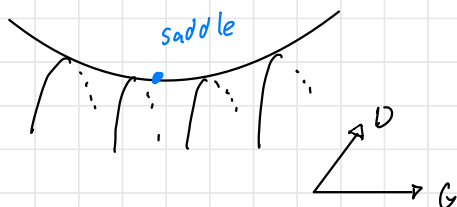
$= \frac{1}{n} \sum_{i=1}^n \log D(x_i) + \frac{1}{n} \sum_{i=1}^n \log (1 - D(x_i^*))$

↓
 $G(\tilde{z}_i)$

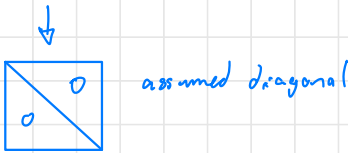
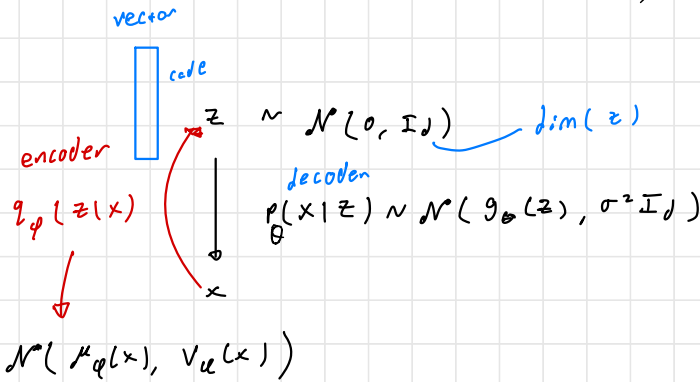
• Value:

$$V(G, D) = \frac{1}{n} \sum_{i=1}^n \log P(x_i) + \frac{1}{n} \sum_{i=1}^n \log (1 - D(G(\tilde{z}_i)))$$

min_G max_D V Adversarial Game



• VAE (Variational auto-encoder)



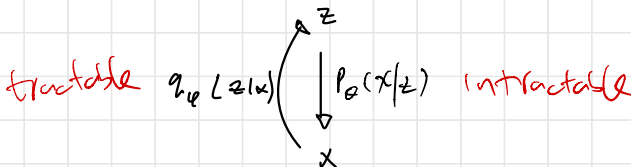
So $z = \mu_\phi(x) + \sqrt{V_\phi(x)} \cdot \underbrace{e}_{\mathcal{N}(0, I_D)}$

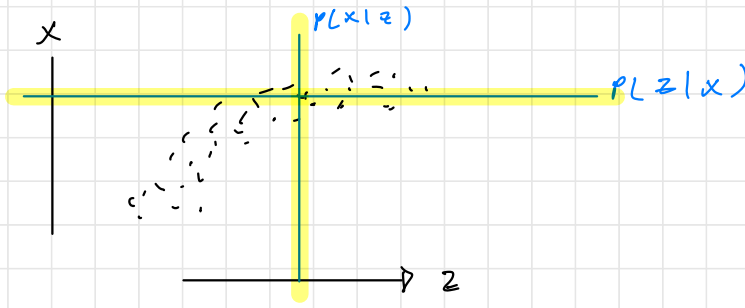
$p_\theta(x, z) = p(z) p_\theta(x|z)$ factorization

$p_\theta(x) = \int p_\theta(x, z) dz$ marginalization (intractable)

$p_\theta(z|x) = \frac{p_\theta(x, z)}{p_\theta(x)}$ conditioning (intractable)

$q_\mu(z|x)$ serves as tractable approx to $p_\theta(z|x)$





• Maximum Likelihood

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_{\theta}(x_i)$$

intractable

density estimation

• Translation:

- If $n \rightarrow \infty$

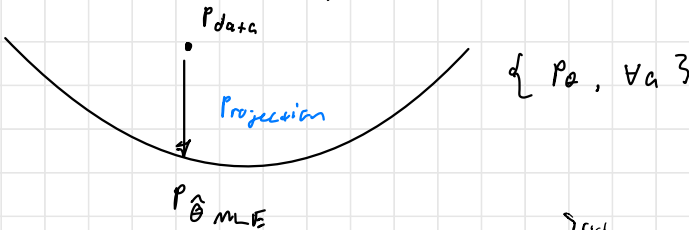
$$L(\theta) \doteq \mathbb{E}_{p_{\text{data}}(x)} [\log p_{\theta}(x)]$$

$x_1 \dots x_n \dots x_n \stackrel{i.i.d.}{\sim} p_{\text{data}}(x)$

$$\max_{\theta} L(\theta) = \min_{\theta} \text{const} - L(\theta)$$

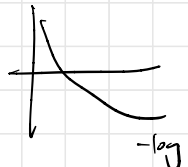
$$\mathbb{E}_{p_{\text{data}}} [\log p_{\text{data}}(x) - \log p_{\theta}(x)]$$

$$= \mathbb{E}_{p_{\text{data}}} \left[\log \frac{p_{\text{data}}(x)}{p_{\theta}(x)} \right] = D_{KL}(p_{\text{data}}(x) \parallel p_{\theta}(x))$$



$$D_{KL}(p \parallel q) = \mathbb{E}_p \left[\log \frac{p(x)}{q(x)} \right] = \mathbb{E}_p \left[-\log \frac{q(x)}{p(x)} \right] \geq -\log \left(\mathbb{E}_p \left[\frac{q(x)}{p(x)} \right] \right) = 0$$

$$\int \frac{q(x)}{p(x)} p(x) dx = \int q(x) dx = 1$$



• MLE :

$$\min_{\theta} D_{KL} (P_{data}(x) | P_{\theta}(x))$$

Intractable

VAE :

$$\min_{\theta, q} D_{KL} (P_{data}(x, z) | P_{\theta}(x, z))$$

$p(z) P_{\theta}(x|z)$
decoder

Tractable

$$P_{data}(x) \cdot q_{\phi}(z|x)$$

encoder

tractable

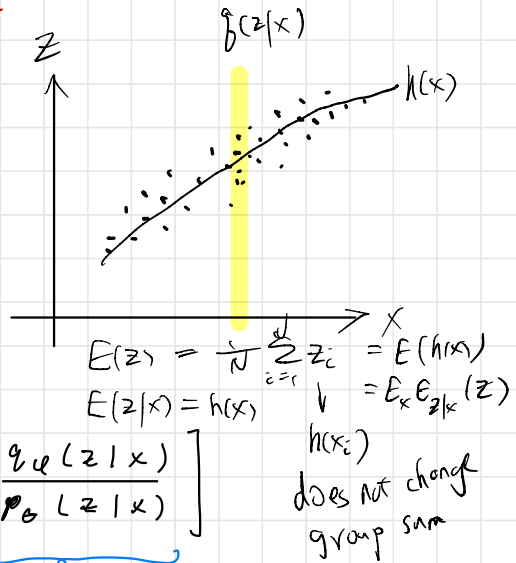
• Understanding :

$$D_{KL} (P_{data}(x, z) | P_{\theta}(x, z))$$

$$= \mathbb{E}_{P_{data}} \left[\log \frac{P_{data}(x, z)}{P_{\theta}(x, z)} \right]$$

$$= \mathbb{E}_{P_{data}} \left[\underbrace{\frac{\log P_{data}(x)}{P_{\theta}(x)}}_{MLE} + \log \frac{q_{\phi}(z|x)}{P_{\theta}(z|x)} \right]$$

Expect



$$= \mathbb{E}_{P_{data}} \left[\log P_{data}(x) - \log P_{\theta}(x) + \log \frac{q_{\phi}(z|x)}{P_{\theta}(z|x)} \right]$$

$$= \mathbb{E}_{P_{data}(x)} \mathbb{E}_{q_{\phi}(z|x)} \left[\right]$$

$$= \mathbb{E}_{P_{data}(x)} \left[\log P_{data}(x) - \log P_{\theta}(x) + \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{q_{\phi}(z|x)}{P_{\theta}(z|x)} \right] \right]$$

$$= \mathbb{E}_{p_{\text{data}}(x)} [\log p_{\text{data}}(x) - \left(\log p_{\theta}(x) - D_{\text{KL}}(q_{\phi}(z|x) | p_{\theta}(z(x))) \right)]$$

Intractable

- Approx $\log p_{\theta}(x) \approx h_T$

$$\text{ELBO: } \log p_{\theta}(x) - D_{\text{KL}}(q_{\phi}(z|x) | p_{\theta}(z(x)))$$

Tractable

Evidence Lower Bound

Problem: (1) $p_{\text{data}}(x) q_{\phi}(z|x) \rightarrow q(z)$ not $N(0, I)$

(2) $p_{\theta}(x(z)) \sim N(g_{\theta}(x), \sigma^2 I_{\theta})$ not accurate
 ↓
 big

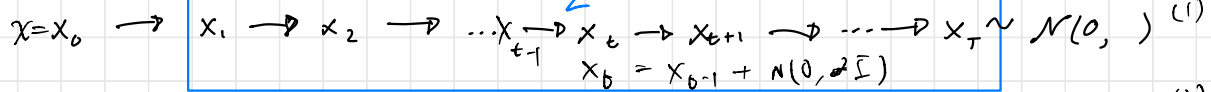
• Diffusion: more careful VAE

no ϕ to learn

Encoder

$q(z|x_0)$ add noise

- clean image



Supervised learning $P(x_{t+1} | x_t) \sim N(x_t + \sigma \nabla_x \log p_{\theta}(x_t), \sigma^2 I)$ accurate (2)

Decoder: $p(z) p(x|z)$

When noise is small