

According to the graphical model  $p(X_{1:T}|Z_{1:T})$  as  $\prod_{t=1}^T p(X_t|Z_{1:t})$ . Apply d-separation:  $p(X_1, X_2, X_3, X_4 | Z_1, Z_2, Z_3, Z_4) = p(X_1|Z_1)p(X_2|X_1, Z_2)p(X_3|X_2, Z_3)p(X_4|X_3, Z_4)$

$$\hookrightarrow q(X_{1:T}|Z_{1:T}) \text{ as } \prod_{t=1}^T q(Z_t|Z_{1:t-1}) \Rightarrow q(Z_1, Z_2, Z_3, Z_4 | X_1, X_2, X_3, X_4) = q(Z_1|X_1) \cdot q(Z_2|X_1, Z_1) \cdot q(Z_3|X_2, Z_1, Z_2) \cdot q(Z_4|X_3, Z_2, Z_3, Z_4) = \prod_{t=1}^T q(Z_t|X_{1:t-1}, Z_{1:t-1})$$

(a) factorize  $p(Z_{1:T}|X_{1:T}) \Rightarrow p(Z_1|X_1) \cdot p(Z_2|Z_1, X_1) \cdot p(Z_3|Z_1, Z_2, X_3) \cdot p(Z_4|Z_1, Z_2, Z_3, X_4)$  (b) joint distribution  
Describe how to evaluate the probability  $p(X_1, X_2, X_3)$  or ...  $p(X_1, X_2, X_3) = p(Z_1) \left| \frac{\partial Z_1}{\partial X_1} \right| p(Z_2|Z_1) \left| \frac{\partial Z_2}{\partial X_1} \right| p(Z_3|Z_1, Z_2) \left| \frac{\partial Z_3}{\partial X_1} \right| p(Z_4|Z_1, Z_2, Z_3) \left| \frac{\partial Z_4}{\partial X_1} \right|$

How would train the flow? maximum log likelihood  $L = \text{Ex}[\log p(x)]$  or find  $\text{argmin}_{\theta} \text{Ex}[-\log p_\theta(x)] = \text{Ex}[-\log p(\theta|x)] - \log \det \left| \frac{\partial f_\theta(x)}{\partial x} \right|$

Sample 1: train 2i using flow model  $\Rightarrow$  use trained inverse function  $f^{-1}$  to generate  $x_i$  from  $z_i$ :  $x_1 = f_1^{-1}(z_1)$ ,  $x_2 = f_2^{-1}(z_2)$ ,  $x_3 = f_3^{-1}(z_3)$

Design an invertible mapping from  $X_i \rightarrow Z_i$ : flow model use invertible  $f_1 \rightarrow z_1 = f_1(x_1)$ ,  $f_2 \rightarrow z_2 = f_2(x_2, z_1)$ ,  $f_3 \rightarrow z_3 = f_3(x_3, z_2, z_1)$

Design an encoding distribution  $q(Z_{1:T}|X_{1:T})$  to approx  $p(Z_{1:T}|X_{1:T})$ ,  $\hat{z}_t$  based on  $X_{1:t}, Z_{1:t}$ , factorize of  $q(Z_{1:T}|X_{1:T}) = \prod_{t=1}^T q(Z_t|X_{1:t}, Z_{1:t}, X_{1:t})$

Train this latent model by maximum ELBO. Describe network CNN for encoding, decoding for prior distribution. Training object:  $t \mapsto$   
(c) Encoder = CNN  $X_{1:T} \rightarrow Z_{1:T}$ , Decoder  $Z_{1:T} \rightarrow X_{1:T}$  objective = Maximum ELBO, prior distribution = Gaussian distribution, ELBO:  $E_{q(Z_{1:T}|X_{1:T})} [\log p(X_{1:T}|Z_{1:T})] - \text{KL}(q(Z_{1:T}|X_{1:T}) || p(X_{1:T}))$

dropout = 隨機棄用一些 hidden unit 不參與計算。是一種 regularization 手段

gradient vanish problem: (cause ① activation 邊坡太小, ex: sigmoid 在極端值的導數  $\rightarrow 0$  solve ② 使用 ReLU activation, batch normalization, Residual Network  
③ 神經路太深, 前面單路梯度極小

CNN contribution: ① Local Connectivity = CNN 利用局部連接來捕捉局部的特徵

② Parameter Sharing: 在各層中權重在不同的空間位置共享, 減少參數量, 提高效率

pooling = 用平均法減少 feature graph 的維度, 從而降低

計算複雜度, 同時保持重要特徵

Training VAE  
try to maximize a ELBO. Explain main idea, variational lower bound:  $L = E_{z \sim g(z|x; \theta)} [\log p(x; \theta) - \text{KL}(g(z|x; \theta) || p(z))]$

What distribution approximate  $q(z|x)$ ? Gaussian distribution, this is an assumption.

Explain the notion of reparam trick: 將隨機變量  $z \sim g(z|\mu)$  重新寫成一個確定性變量的函數  $z = g(\varepsilon, x)$ ,  $\varepsilon$  是從高斯分佈抽樣

Posterior  $q(z|x) = p(z)$ ? ans: False, when maximum variational lower bound,  $g(z|x)$  shouldn't be identity to  $p(z)$ . However it should catch the posterior distribution

How would you evaluate KL divergence if prior  $p(z)$  is replace with Gaussian Mix...

$$\text{KL}(q(z|x) || p(z)) = E_{q(z|x)} \left[ \log \frac{q(z|x)}{p(z)} \right] \text{ substitute } p(z) \text{ with GMN}$$

	VAE	GAN	WGAN	Flow	Diffusion	AE
stochastic	maximum Variational lower bound	minimax loss	minimize Earth Mover's Distance	maximum log likelihood	minimize loss that combine reconstruction at each step	minimize reconstruction loss
stochastic						
stochastic						
deterministic						
stochastic						
deterministic						

$\checkmark \quad \times \quad \text{Reconstruction loss, KL Divergence}$

$\times \quad \times \quad \text{Fréchet Inception Distance (FID)}$

$\times \quad \times \quad \text{Wasserstein distance}$

$\checkmark \quad \checkmark \quad \text{Log likelihood}$

$\checkmark \quad \checkmark \quad \text{FID}$

$\times \quad \times \quad \text{Reconstruction loss}$

$\checkmark \quad \text{generating latent code of specific data sample}$

Evaluate: Probability of specific data sample

In evaluating KL divergence between ground-truth  $p(z)$  and learned distribution  $q(z)$  ... while  $p(z)$  has 2 peak.

$$\text{KL}(p||q) = - \int p(x) \log \left[ \frac{q(x)}{p(x)} \right] dx \geq E_{x \sim p} [\log p(x) - \log q(x)]$$

$$\text{① } \min(\text{KL}(p||q)) = \min_{z \sim p} (E[\log p(z) - \log q(z)]) = \min_{z \sim p} (E[\log \frac{p(z)}{q(z)}]), \text{ KL Divergence} \geq 0$$

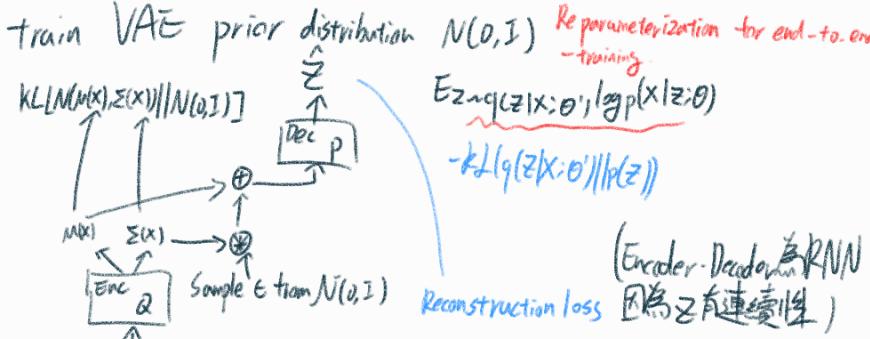


$\frac{p(z)}{q(z; \theta)}$  越接近 1, KL 越接近 0

$$\text{② } \min(\text{KL}(q||p)) = \min_{z \sim q} (E[\log q(z) - \log p(z)]) = \min_{z \sim q} (E[\log \frac{q(z)}{p(z)}]). \text{ q 做 sample 是能局部接近 p}$$



沒辦法知道 p 的整体



ReLU  $\begin{cases} \text{差異 } x < 0 \\ \text{輸出 } 0, \text{無梯度} \end{cases}$

LeakyReLU  $\begin{cases} \text{一條很小斜率的直線} \end{cases}$

ELU  $\begin{cases} \text{幅度小的曲線} \end{cases}$

Sigmoid  $\begin{cases} \text{一幅圖} \end{cases}$

Conv feature map  $(K_h, K_w)$   
stride  $(S_h, S_w)$   
padding  $(P_h, P_w)$

$$H_{\text{out}} = \left\lceil \frac{H_{\text{in}} - K_h + 2 \times P_h}{S_h} \right\rceil + 1$$

$$W_{\text{out}} = \left\lceil \frac{W_{\text{in}} - K_w + 2 \times P_w}{S_w} \right\rceil + 1$$

$$\text{parameters} = K_h \cdot K_w \cdot \text{Input}_c \cdot \text{Output}_c$$

Compare RNN-based, VAE, GAN, flow model, diffusion model

allow exact evaluation of likelihood = VAE、Flow-based-RNN

training obj RNN: 最大化下個時間步驟的條件概率

VAE: 最大化 ELBO

GAN: 最小化生成器和discriminator 的对抗损失 最大化 JS-divergence

Flow-based: maximum log likelihood

Diffusion: 最小化每一步 Diffusion 的重構誤差 通常也包含KL-divergence

in terms of generation process, describe how to draw sample

RNN-based auto regressive = 依賴先前 timestep 的 data 來成下一個 timestep 的 data

VAE = 在 latent space 中采樣，生成真實 data 的圖像

Flow-based = 通過可逆地 latent space 樣本轉為 data 樣本，這個過程是確定性的

Diffusion = 逐步降噪生成圖像

major advantages of diffusion as compared to GAN

Diffusion 可以生成更高質量的圖像，訓練不依賴於 generator & discriminator 更穩定

過大的 beta 會使每一步的 noise 變化過大，導致逆向難以準確重構

GAN vs WGAN

Critic

GAN → Maximum JS-divergence

相似

WGAN → Minimize Wasserstein Distance

Transformer vs. RNN

Transformer 是一種基於注意力機制的網路架構，特別適用於處理序列資料

RNN (RNN)

Transformer 可以並行處理，RNN 要逐個 timestep 處理

應用例子：Transformer 在 NLP 任務中非常成功，如 GPT

Dropout, vs L<sub>1</sub>/L<sub>2</sub> regularization: L<sub>1</sub>, L<sub>2</sub> 正則化對 loss function 做修改，Dropout 直接對網路做隨機忽略

應用例子：CNN 中 Fully Connect 之後使用 Dropout 減少 overfitting.

padding: 在輸入的圖邊界添加 extra pixel 以控制輸出 feature map 大小，for instance: Kernel Size (3,3) Img\_sz (5,5)  
Input sz=Output sz Some padding Output\_sz (5,5)  
No padding. No padding Output\_sz (3,3)

attention

↳ Self-attention 計算輸入的元素和其他元素之間的關係來重新加權。

↳ Multi-head-attention 將 Self-attention 利用多次，並縮短距離，以得到更豐富的特徵

例子：NLP 中 attention 可以使 model專注於句子中當前單字最相關的部分

对比

CNN 通常用於圖像處理，但 attention 也能用於捕捉不同區域間的依賴性

RNN: attention 有利於捕捉更長距離依賴的關係