

Generative: AR

$$P_\theta(x) = \prod_{t=1}^T P_\theta(x_t | x_{1:t-1})$$

$$= P_\theta(x_2 | x_{1:1}) P_\theta(x_3 | x_{1:2}) P_\theta(x_4 | x_{1:3}) \cdots P_\theta(x_T | x_{1:T-1})$$

$\bullet \bullet$                        $\bullet \bullet \bullet$                        $\bullet \bullet \bullet \bullet$

$$\log P_\theta(x) = \sum_{t=1}^T \log P_\theta(x_t | x_{1:t-1})$$

Generative: flow

$$x = g(z), \quad z = f(x)$$

$$P(x) = \tilde{P}(f(x)) \left| \det Dg(f(x)) \right|^T = \tilde{P}(f(x)) \left| \det Df(x) \right|$$

$$Dg(z) = \frac{\partial g}{\partial z}$$

$$\det Df(x) = \prod_{i=1}^N \det Df_i(x_i)$$

$\hookrightarrow x_i = g_i \circ \cdots \circ g_1(z)$   
 $= f_{i+1} \circ \cdots \circ f_N(x)$

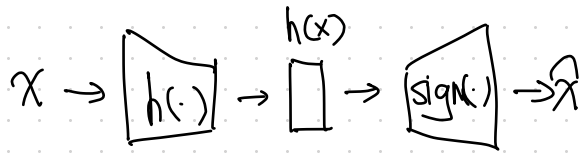
$$z \rightarrow g_1 \rightarrow g_2 \rightarrow \cdots \rightarrow g_i \rightarrow x_i \leftarrow f_N \leftarrow f_{N-1} \leftarrow \cdots \leftarrow f_{i+1} \leftarrow x$$

$$z \rightarrow g_1 \rightarrow g_2 \rightarrow \cdots \rightarrow g_N \rightarrow x$$

$$z \leftarrow f_N \leftarrow f_{N-1} \leftarrow \cdots \leftarrow f_1 \leftarrow x$$

Generative: AE

$$x \rightarrow h(\cdot) \rightarrow \text{sigm}(\cdot) \rightarrow \hat{x}$$

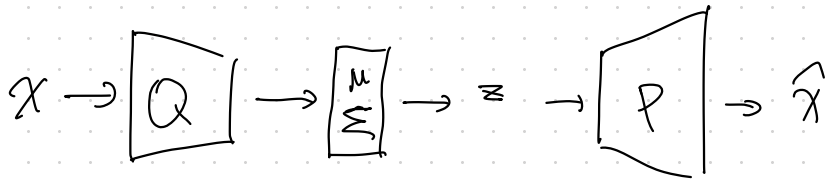


$$h(x) = g(b + W \cdot x) \quad \hat{x} = \text{sigm}(c + V \cdot h(x))$$

Loss

- e.g. binary, ~~cross~~ entropy

Generative : VAE



$$\begin{aligned} \ln p_{\theta}(x) &= \ln \int p_{\theta}(x|z) p(z) dz \\ &= \ln \int \frac{q_{\phi}(z|x)}{q_{\phi}(z|x)} p_{\theta}(x|z) p(z) dz \\ &\geq \underbrace{-\left[ D_{KL}(q_{\phi}(z|x) || p(z)) - \mathbb{E}_{q_{\phi}}(\ln p_{\theta}(x|z)) \right]}_{\substack{F(x) \leftarrow \text{free energy} \\ -F(x) \leftarrow \text{ELBO}}} \end{aligned}$$

← reconstruction error

$\theta$  : model for inference

$\phi$  : variational approximation

Reparametrization trick:

reparametrization of latent var

e.g.,

$$q_{\phi}(z) = \mathcal{N}(z|\mu, \sigma^2), \quad \phi = \{\mu, \sigma^2\}$$

# VAE loss

$$-\ln P(x) + D[q(z|x) \| P(z|x)] = -E_{z \sim q} [\ln P(x|z)] + D[q(z|x) \| P(z)]$$

↓  
bits for  
constructing  $x$   
using ideal  
coding

penalty

as  $q$  and  $p$   
are not  
necessarily  
the same,  
i.e.,  $q$  is  
only sub-optimal

information to  
reconstruct  $x$   
from  $z$   
using ideal coding

construct  $z$ .

extra information  
about  $x$  if we  
get  $x$  using  $z$  from  
 $p(z|x)$  instead of  
 $p(z)$

# Contrastive

learn to compare

NCE : Noise Contrastive Estimation

$$\mathcal{L} = E_{x, x^+, x^-} \left[ -\log \frac{C(x, x^+)}{C(x, x^+) + C(x, x^-)} \right]$$

$x^+$ : similar to  $x$

$x^-$ : dissimilar to  $x$

e.g. 
$$\mathcal{L} = E_{x, x^+, x^-} \left[ -\log \left( \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$f$ : encoder

# Contrastive: Context Instance

## Predict Spatial Relation

- jigsaw
- angle of image

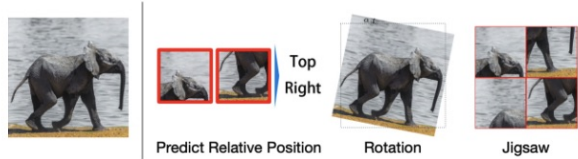


Fig. 8: Three typical methods for spatial relation contrast: predict relative position [37], rotation [43] and solve jigsaw [67], [87], [92], [141].

## Max Mutual Info

$$MI: I(X;Y) = E_{p_{X,Y}} \log \frac{p_{X,Y}}{p_X p_Y}$$

$$I \sim D_{KL}(P_{X,Y} \| P_X P_Y)$$

$$I(X;Y) = H(X) - H(X|Y)$$

reduction of uncertainty of X after observing Y

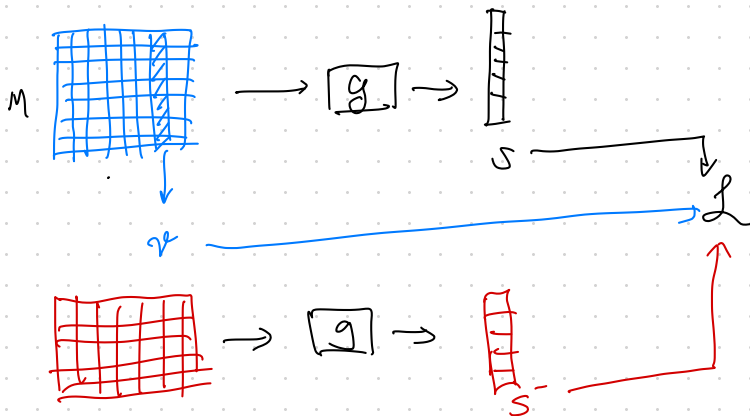
Max MI models:

$$\max_{g_1 \in \mathcal{G}_1, g_2 \in \mathcal{G}_2} I(g_1(x_1), g_2(x_2))$$

# Max MI : DeepInfoMax

$f(x) \in \mathbb{R}^{M \times M \times d}$   
 $\uparrow$   
 image

a feature vector  $u \in \mathbb{R}^d$  from  $f(x)$   
 encoding of img



$f(x) \Rightarrow [g] \xrightarrow{s} \text{context of } v$

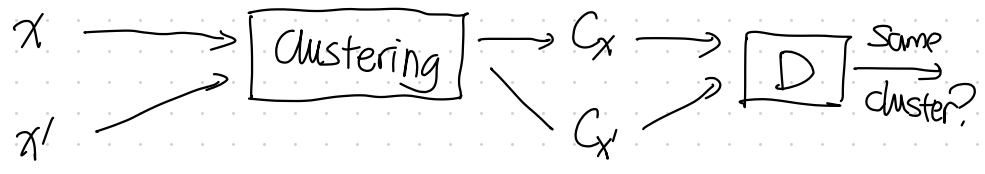
$f(x') \Rightarrow [g] \xrightarrow{s^-} \text{context of } v^-$

$$\mathcal{L} = \mathbb{E}_{x,x'} \left[ -\log \frac{e^{v^T \cdot s}}{e^{v^T \cdot s} + e^{v^T \cdot s^-}} \right]$$

Contrastive : Instance Instance

## Cluster Discrimination

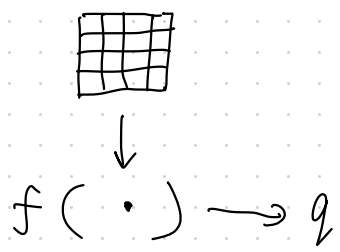
Deep Cluster



## Instance Discrimination

MoCo

Encoder  $q = f_q(x)$



loss:

$$k_+ = f_k(x) \quad k_i = f_k(x_i)$$

$$\mathcal{L} = -\log \frac{\exp(q \cdot \frac{k_+}{\tau})}{\sum_{i=0}^K \exp(q \cdot \frac{k_i}{\tau})}$$

*all negative samples*

two encoders

- query :  $\rightarrow q$
- key :  $\rightarrow k$ , queue of data



# SimCLR

mini-batch ( $N$ )

↓ augment data  
to  $2N$

$$N \begin{pmatrix} \vdots & \vdots \\ \hat{x}_i & \hat{x}_r \\ \vdots & \vdots \end{pmatrix} \begin{matrix} \} \text{negative} \\ \leftarrow \text{positive} \\ \} \text{negative} \end{matrix}$$

$$\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [l_{2i-1, 2i}, l_{2i, 2i-c}]$$

$$l_{i,i'} = -\log \frac{\exp(\text{sim}(\hat{x}_i, \hat{x}_{i'})/\tau)}{\sum_{\substack{k=1 \\ k \neq i}}^N \exp(\text{sim}(\hat{x}_i, \hat{x}_k)/\tau)}$$

Adversarial : Complete Input

GAN

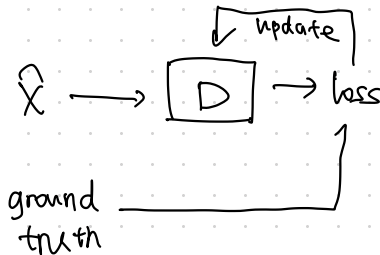
task:  $\{ , Y \rightarrow X$   
noise class feature

min max game : two players  $G, D$ ;  $\min_G \max_D V(D, G)$  } worst case  $G$  min  $V$   
 $G$ : fool  $D$  } then find  $D$  that max  $V$   
 $D$ : min discrimination error

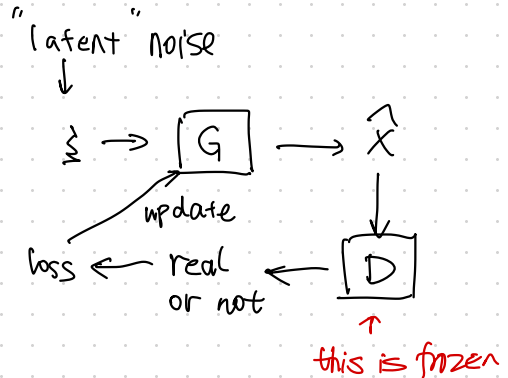
$$\min_G \max_D E_{x \sim P_{\text{data}}} [\log D(x)] + E_{z \sim P_z} [\log (1 - D(G(z)))]$$

Alternating training

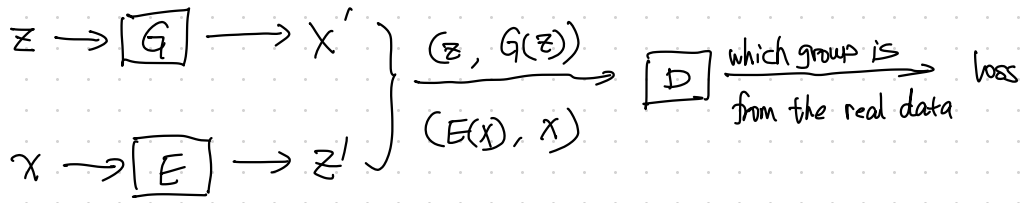
Train discriminator



Train Generator



# BiGAN & ALI



$$\Rightarrow E \rightarrow G^{-1}$$