



电子科技大学  
University of Electronic Science and Technology of China



# Information Theory in Deep Neural Networks

Wei Han



Data Mining Lab,  
Big Data Research Center, UESTC  
Email: [weihan@std.uestc.edu.cn](mailto:weihan@std.uestc.edu.cn)

- Preliminary of Information Theory
- Information between Data and Representation
- Information Bottleneck & Mutual Information
- Disentangled Representation via Mutual Information
- Generalization Bound from Mutual Information



# Preliminary

- Uncertainty



- Shannon Entropy

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x)$$

- Uncertainty



$$H(X) = -\frac{1}{2} \times \log \frac{1}{2} - \frac{1}{2} \times \log \frac{1}{2} = 1$$

$$P\{Y = 1, 2, \dots, 6\} = \frac{1}{6} \quad H(Y) = 6 \times \left(-\frac{1}{6} \log \frac{1}{6}\right) = \log 6$$

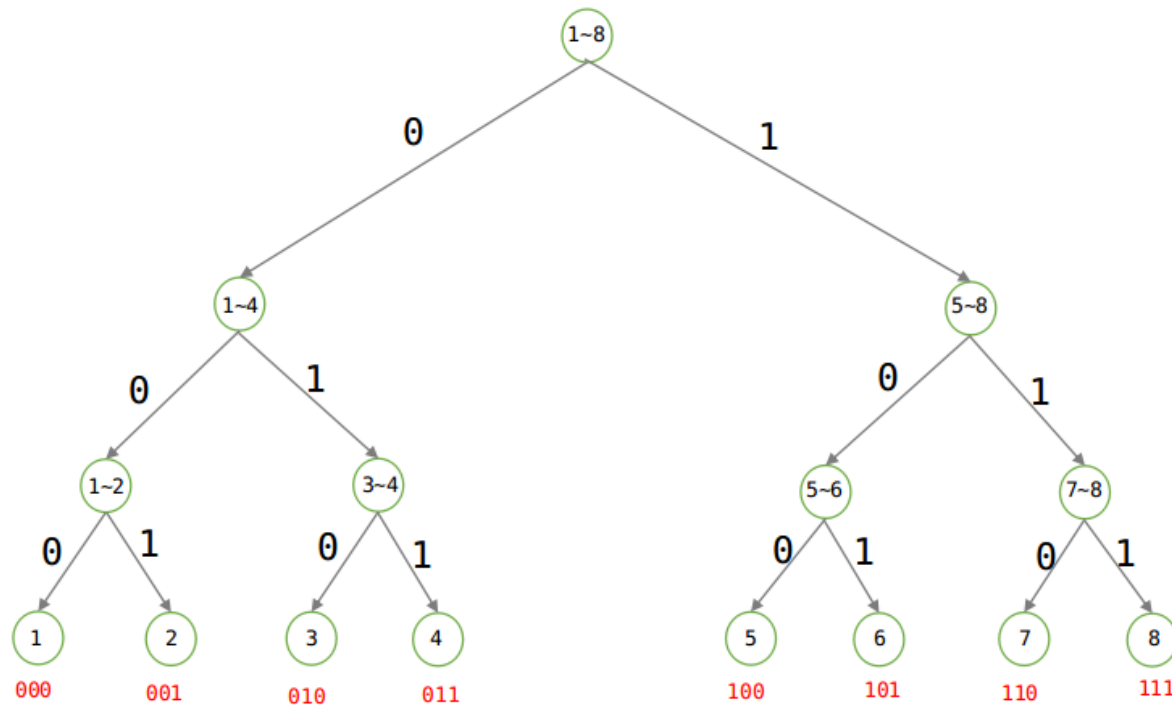
$$H(X) = 1 = \log 2 < H(Y) = \log 6$$

# 1.2. Mean Coding Length



- Minimum MCL

$$H(x) = - \sum_{x \in X} p(x) \log p(x) = E_{x \sim P}[-\log p(X)]$$



- Coding  $P$  with  $Q$

$$H(P, Q) = E_{x \sim P}[-\log Q(X)] = - \sum_{x \in X} p(x) \log q(x)$$

Coding

Minimize MCL

- Distribution Gap

$$\begin{aligned} KL(P||Q) &= \sum P(x) \log \frac{P(x)}{Q(x)} = \sum_{x \in X} p(x) \log p(x) - \sum_{x \in X} p(x) \log q(x) \\ &= \underbrace{-H(P)} + H(P, Q) \end{aligned}$$



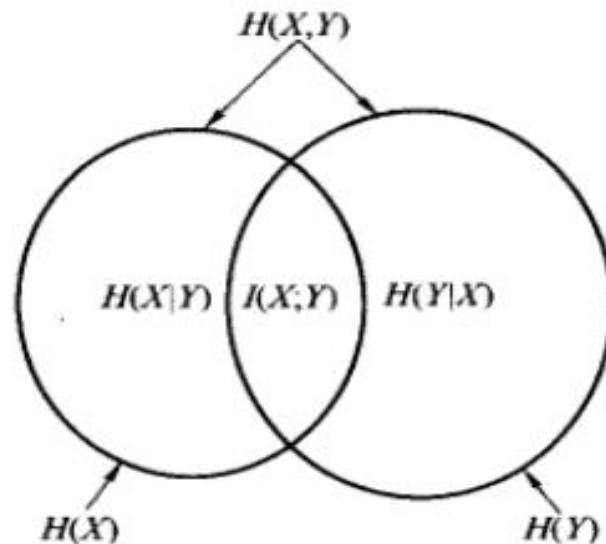
Data are given



- Dependence

$$I(X; Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x, y) \log \left( \frac{p(x, y)}{p(x) p(y)} \right)$$

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





# Information between Data and Representation

- Information Entropy

$$H_c = - \sum_{c \in \text{汉字}} p_c \log p_c$$

Chinese character *9.65bit* vs. English character *4.03bit*

From characters to words

$$- \sum_{w \in \text{语料}} \log p_w$$

- Information Entropy

$$-\sum_{w \in \text{语料}} \log p_w < -\sum_{c \in \text{语料}} \log p_c \longrightarrow -\sum_{w \in \text{词表}} N_w \log p_w < -\sum_{c \in \text{字表}} N_c \log p_c$$

$$\begin{aligned} \mathcal{L} &= -\sum_{w \in \text{词表}} \frac{N_w}{\text{总字数}} \log p_w \\ &= \left( -\sum_{w \in \text{词表}} \frac{N_w}{\text{总词数}} \log p_w \right) \div \left( \frac{\text{总字数}}{\text{总词数}} \right) \\ &= \left( -\sum_{w \in \text{词表}} \frac{N_w}{\text{总词数}} \log p_w \right) \div \left( \frac{\sum_{w \in \text{词表}} N_w l_w}{\text{总词数}} \right) \\ &= \left( -\sum_{w \in \text{词表}} \frac{N_w}{\text{总词数}} \log p_w \right) \div \left( \sum_{w \in \text{词表}} \frac{N_w}{\text{总词数}} l_w \right) \\ &= \frac{H}{l} \end{aligned}$$

$$-\sum_{c \in \text{字表}} \frac{N_c}{\text{总字数}} \log p_c = -\sum_{c \in \text{字表}} p_c \log p_c$$

- Localize / Segment

$$\mathcal{H} = - \sum_{w \in \text{词表}} p_w \log p_w, \quad l = \sum_{w \in \text{词表}} p_w l_w$$
$$\mathcal{L} = \frac{\mathcal{H}}{l} = \frac{- \sum_i p_i \log p_i}{\sum_i p_i l_i}$$

To combine character  $a$  and  $b$  as a word

$$\tilde{p}_{ab} = \frac{N_{ab}}{\tilde{N}} = \frac{p_{ab}}{1 - p_{ab}}$$
$$\tilde{p}_a = \frac{\tilde{N}_a}{\tilde{N}} = \frac{p_a - p_{ab}}{1 - p_{ab}}, \quad \tilde{p}_b = \frac{\tilde{N}_b}{\tilde{N}} = \frac{p_b - p_{ab}}{1 - p_{ab}}$$
$$\tilde{p}_i = \frac{N_i}{\tilde{N}} = \frac{p_i}{1 - p_{ab}}, \quad (i \neq a, b)$$

- Minimize Entropy

$$\begin{aligned}\tilde{H} &= -\frac{1}{1-p_{ab}} \left\{ p_{ab} \log\left(\frac{p_{ab}}{1-p_{ab}}\right) + \right. \\ &\quad \left. \sum_{i=a,b} (p_i - p_{ab}) \log\left(\frac{p_i - p_{ab}}{1-p_{ab}}\right) + \sum_{i \neq a,b} p_i \log\left(\frac{p_i}{1-p_{ab}}\right) \right\} \\ &= \frac{1}{1-p_{ab}} (\mathcal{H} - \mathcal{F}_{ab}) \quad \begin{array}{c} \uparrow \\ \downarrow \end{array}\end{aligned}$$

$$\begin{aligned}\mathcal{F}_{ab} &= p_{ab} \log \frac{p_{ab}}{p_a p_b} - (1-p_{ab}) \log(1-p_{ab}) \\ &\quad + \sum_{i=a,b} (p_i - p_{ab}) \log\left(1 - \frac{p_{ab}}{p_i}\right)\end{aligned}$$

- Minimize Entropy

$$\begin{aligned}\tilde{l} &= \frac{p_{ab}}{1 - p_{ab}}(l_a + l_b) + \sum_{i=a,b} \frac{p_i - p_{ab}}{1 - p_{ab}} l_i + \sum_{i \neq a,b} \frac{p_i}{1 - p_{ab}} l_i \\ &= \frac{l}{1 - p_{ab}}\end{aligned}$$

After knowing the change, the change of  $L$ :

$$\frac{\tilde{\mathcal{H}}}{\tilde{l}} - \frac{\mathcal{H}}{l} = -\frac{\mathcal{F}_{ab}}{l} \uparrow \quad \uparrow$$

- Minimize Entropy

$$\mathcal{F}_{ab} \approx \mathcal{F}_{ab}^* = p_{ab} \left( \ln \frac{p_{ab}}{p_a p_b} - 1 \right)$$

$$\mathcal{F}_a \approx \mathcal{F}_a^* = p_a \left( \ln \frac{p_a}{\prod_{i \in a} p_i} - 1 \right)$$

$$PMI(a, b) = \ln \frac{p_{ab}}{p_a p_b} > 1$$

Segment



$$\ln \frac{p_{ab}}{p_a p_b} < \text{min\_pmi}$$

Split



## 2.2. Minimum Entropy in Representation



- Library Model

$$S = \sum_{i,j} p(i) d(i)$$



- What if more than one book?

$$S = \sum_{ij} p(i)p(j|i)[d(i) + d(i,j)] = \sum_{ij} p(i,j)[d(i) + d(i,j)]$$

- book2vec

$$S = \min_{\mathbf{v}} \sum_{ij} p(i)p(j|i) [\|\mathbf{v}_i\| + \|\mathbf{v}_i - \mathbf{v}_j\|] = \sum_{ij} p(i,j) [\|\mathbf{v}_i\| + \|\mathbf{v}_i - \mathbf{v}_j\|]$$

s.t.  $\forall i, j, \|\mathbf{v}_i - \mathbf{v}_j\| \geq d_{\min}$

$$S = \sum_{ij} p(i)p(j|i) \underline{f(\mathbf{v}_i, \mathbf{v}_j)} = \sum_{ij} p(i,j) \underline{f(\mathbf{v}_i, \mathbf{v}_j)}$$

- What if more than one book?

$$S = \sum_{ij} p(i)p(j|i)f(\mathbf{v}_i, \mathbf{v}_j) = \sum_{ij} p(i, j)f(\mathbf{v}_i, \mathbf{v}_j)$$

- word2vec – skip gram

$$f(\mathbf{v}_i, \mathbf{v}_j) = -\log \frac{e^{-\|\mathbf{v}_i - \mathbf{v}_j\|^2}}{Z_i}, \quad Z_i = \sum_j e^{-\|\mathbf{v}_i - \mathbf{v}_j\|^2}$$

$$f(\mathbf{v}_i, \mathbf{v}_j) = -\log \frac{e^{\langle \mathbf{v}_i, \mathbf{v}_j \rangle}}{Z_i}, \quad Z_i = \sum_j e^{\langle \mathbf{v}_i, \mathbf{v}_j \rangle}$$

$$S = \sum_{ij} p(i)p(j|i)f(\mathbf{v}_i, \mathbf{v}_j) = - \sum_{ij} p(i)p(j|i) \log q(j|i)$$

- (t-)SNE

Let  $P(i) = c$

$$p(\mathbf{x}_j|\mathbf{x}_i) = \frac{e^{-\|\mathbf{x}_i-\mathbf{x}_j\|^2/2\sigma^2}}{\sum_{j \neq i} e^{-\|\mathbf{x}_i-\mathbf{x}_j\|^2/2\sigma^2}}$$

$$S = - \sum_{i \neq j} p(\mathbf{x}_j|\mathbf{x}_i) \log q(j|i), \quad q(j|i) = \frac{e^{-\|\mathbf{v}_i-\mathbf{v}_j\|^2}}{\sum_{j \neq i} e^{-\|\mathbf{v}_i-\mathbf{v}_j\|^2}}$$

- Generative model from  $Z$  to  $X$

$$p(X) \sim \{X_1, \dots, X_n\}$$

$$p(x) = \int p(x|z)p(z)dz, \quad p(x, z) = \underline{p(x|z)p(z)}$$

$$KL\left(p(x, z) \parallel q(x, z)\right) = \iint p(x, z) \ln \frac{p(x, z)}{q(x, z)} dz dx$$

$$\begin{aligned} KL\left(p(x, z) \parallel q(x, z)\right) &= \int p(x) \left[ \int p(z|x) \ln \frac{p(x, z)}{q(x, z)} dz \right] dx \\ &= \mathbb{E}_{x \sim p(x)} \left[ \int p(z|x) \ln \frac{p(z|x)p(x)}{q(x, z)} dz \right] \end{aligned}$$

- Joint Probability Matching

$$KL(p(x, z) \parallel q(x, z)) = \mathbb{E}_{x \sim p(x)} \left[ \int p(z|x) \ln \frac{p(z|x)p(x)}{q(x, z)} dz \right]$$

$$\ln \frac{p(z|x)p(x)}{q(x, z)} = \ln \frac{p(z|x)}{q(x, z)} + \ln p(x)$$

$$\begin{aligned} \mathbb{E}_{x \sim p(x)} \left[ \int p(z|x) \ln p(x) dz \right] &= \mathbb{E}_{x \sim p(x)} \left[ \ln p(x) \int p(z|x) dz \right] \\ &= \mathbb{E}_{x \sim p(x)} [\ln p(x)] \end{aligned}$$

$$\mathcal{L} = KL(p(x, z) \parallel q(x, z)) - \text{常数} = \mathbb{E}_{x \sim p(x)} \left[ \int p(z|x) \ln \frac{p(z|x)}{q(x, z)} dz \right]$$



- Losses

$$\begin{aligned}
 \mathcal{L} &= \mathbb{E}_{x \sim p(x)} \left[ \int p(z|x) \ln \frac{p(z|x)}{q(x|z)q(z)} dz \right] \\
 &= \mathbb{E}_{x \sim p(x)} \left[ - \int p(z|x) \ln q(x|z) dz + \int p(z|x) \ln \frac{p(z|x)}{q(z)} dz \right] \\
 &= \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim p(z|x)} \left[ - \ln q(x|z) \right] + \mathbb{E}_{z \sim p(z|x)} \left[ \ln \frac{p(z|x)}{q(z)} \right] \right] \\
 &= \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim p(z|x)} \left[ \underbrace{- \ln q(x|z)}_{\text{Loss for Generation}} + \underbrace{KL(p(z|x) \parallel q(z))}_{\text{Uniform Gaussian as a Prior}} \right] \right]
 \end{aligned}$$

Loss for Generation

Uniform Gaussian as a Prior

- Components

$$\mathcal{L} = \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim p(z|x)} \left[ \underbrace{-\ln q(x|z)}_{\text{Bernoulli / Gaussian}} \right] + \underbrace{KL\left(p(z|x) \parallel \underbrace{q(z)}_{\text{Standard Gaussian}}\right)}_{\text{Standard Gaussian}} \right]$$

Bernoulli / Gaussian

Standard Gaussian

$$\hat{p}(z|x) = q(z|x) = \frac{q(x|z)q(z)}{q(x)} = \frac{q(x|z)q(z)}{\int q(x|z)q(z)dz}$$



- $P(z|x)$

$$p(z|x) = \frac{1}{\prod_{k=1}^d \sqrt{2\pi\sigma_{(k)}^2(x)}} \exp\left(-\frac{1}{2} \left\| \frac{z - \mu(x)}{\sigma(x)} \right\|^2\right)$$

$$KL\left(p(z|x) \parallel q(z)\right) = \frac{1}{2} \sum_{k=1}^d \left( \underbrace{\mu_{(k)}^2(x)} + \underbrace{\sigma_{(k)}^2(x)} - \ln \sigma_{(k)}^2(x) - 1 \right)$$

- $P(x/z)$  with Bernoulli distribution

$$p(\xi) = \begin{cases} \rho, & \xi = 1; \\ 1 - \rho, & \xi = 0 \end{cases}$$

$$q(x|z) = \prod_{k=1}^D \left( \rho^{(k)}(z) \right)^{x^{(k)}} \left( 1 - \rho^{(k)}(z) \right)^{1-x^{(k)}}$$

$$-\ln q(x|z) = \sum_{k=1}^D \left[ -x^{(k)} \ln \rho^{(k)}(z) - (1 - x^{(k)}) \ln \left( 1 - \rho^{(k)}(z) \right) \right]$$

---

Cross Entropy Loss

- $P(x|z)$  with Gaussian distribution

$$q(x|z) = \frac{1}{\prod_{k=1}^D \sqrt{2\pi\sigma_{(k)}^2(z)}} \exp\left(-\frac{1}{2} \left\| \frac{x - \mu(z)}{\sigma(z)} \right\|^2\right)$$

$$-\ln q(x|z) = \frac{1}{2} \left\| \frac{x - \mu_k(z)}{\sigma(z)} \right\|^2 + \frac{D}{2} \ln 2\pi + \frac{1}{2} \sum_{k=1}^D \ln \sigma_{(k)}^2(z)$$

$$-\ln q(x|z) \sim \frac{1}{2\sigma^2} \underbrace{\left\| x - \mu_k(z) \right\|^2}_{\text{MSE}}$$

MSE



# Information Bottleneck & Mutual Information

- Information Extraction

$$\max[I(Z; Y) - \beta I(X; Z)]$$



# 3.1. Information Bottleneck



- Loss

$$I(X;Z)$$

$$\begin{aligned} & \int \int p(z|x)\tilde{p}(x) \log \frac{p(z|x)}{p(z)} dx dz \\ &= \int \int p(z|x)\tilde{p}(x) \log \frac{p(z|x)}{q(z)} \frac{q(z)}{p(z)} dx dz \\ &= \int \int p(z|x)\tilde{p}(x) \log \frac{p(z|x)}{q(z)} + \int \int p(z|x)\tilde{p}(x) \log \frac{q(z)}{p(z)} dx dz \\ &= \int \int p(z|x)\tilde{p}(x) \log \frac{p(z|x)}{q(z)} + \int p(z) \log \frac{q(z)}{p(z)} dz \\ &= \int \int p(z|x)\tilde{p}(x) \log \frac{p(z|x)}{q(z)} - \int p(z) \log \frac{p(z)}{q(z)} dz \\ &= \int \tilde{p}(x) KL(p(z|x) \parallel q(z)) dx - \underbrace{KL(p(z) \parallel q(z))}_{>0} \\ &< \int \tilde{p}(x) KL(p(z|x) \parallel q(z)) dx \end{aligned}$$

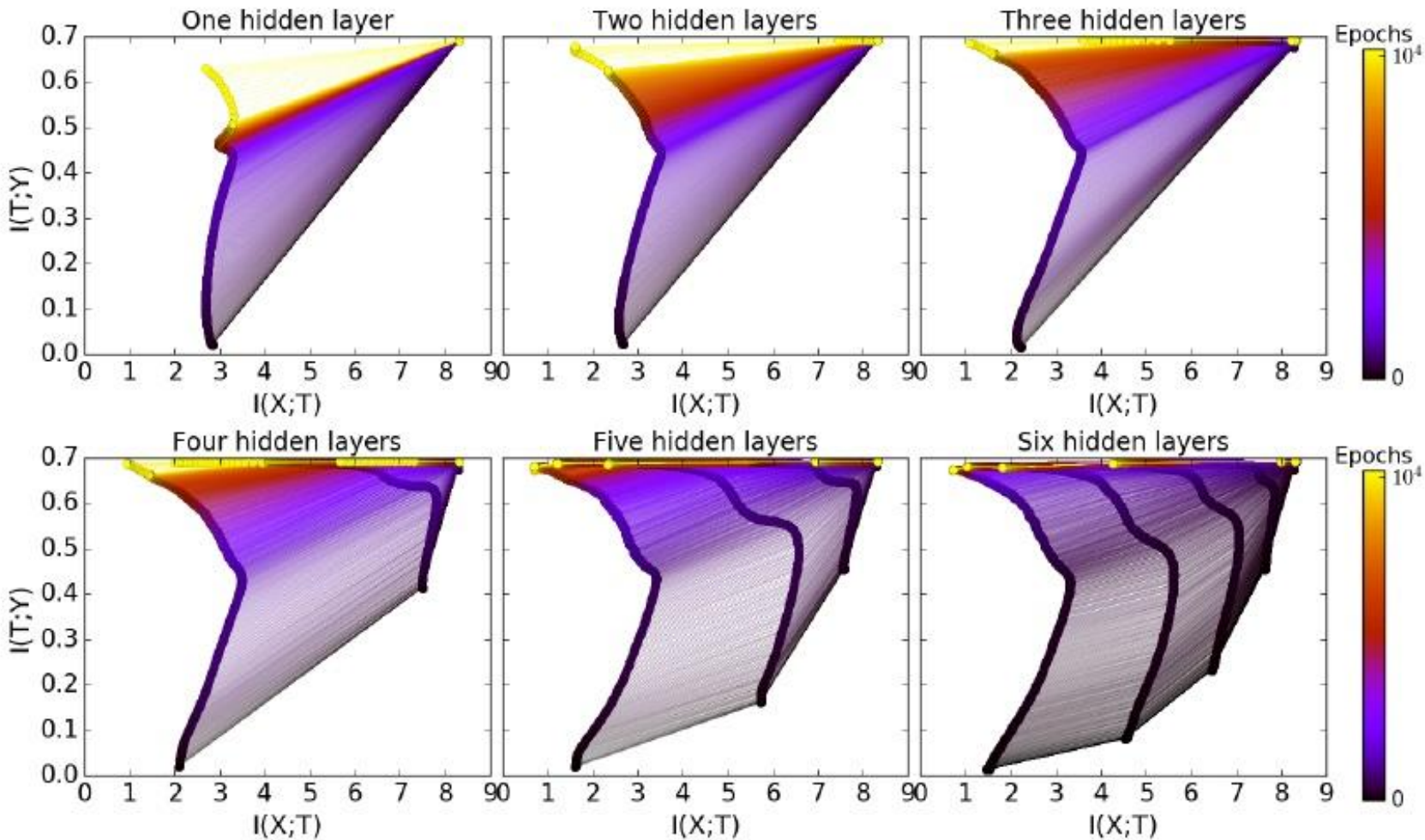
$$KL(p(z|x) \parallel q(z))$$

Standard Gaussian

# 3.1. Information Bottleneck



- Two phase of learning





- AutoEncoder – Keep key information



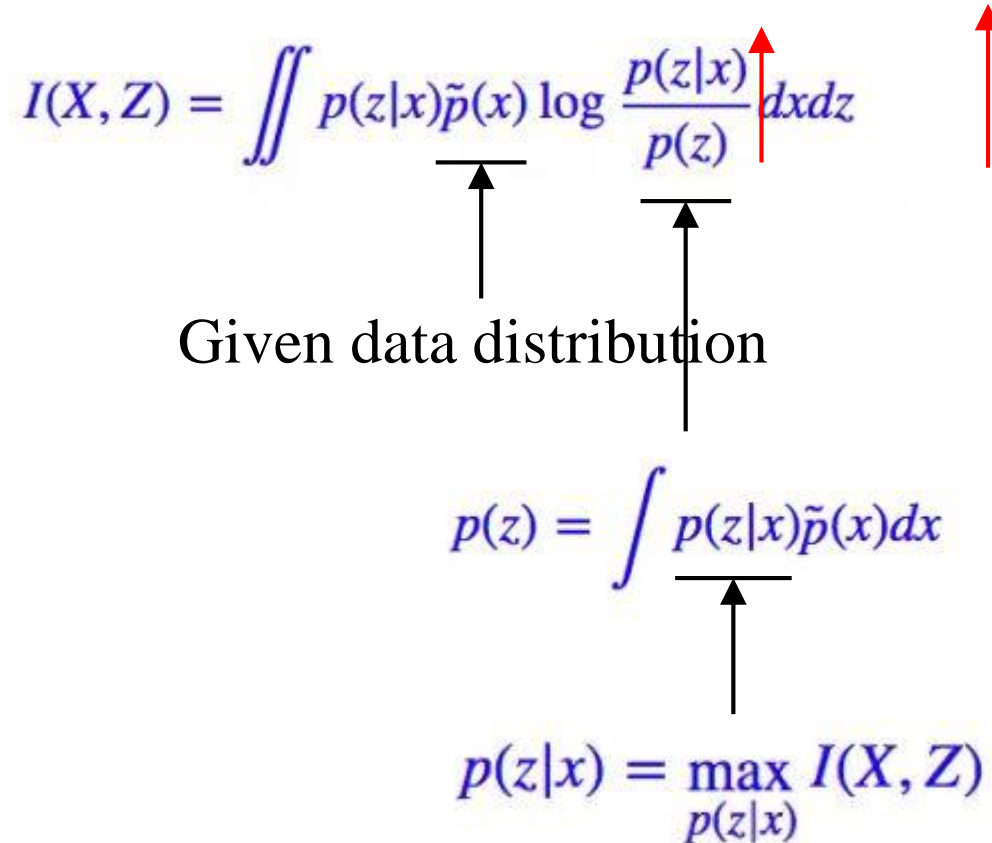
$$p(z|x) = \max_{p(z|x)} I(X, Z)$$



- Maximum Mutual Information

$$I(X, Z) = \iint p(z|x)\tilde{p}(x) \log \frac{p(z|x)}{p(z)} dx dz$$

Given data distribution

$$p(z) = \int p(z|x)\tilde{p}(x) dx$$
$$p(z|x) = \max_{p(z|x)} I(X, Z)$$


- Prior Distribution

$$KL(p(z)||q(z)) = \int p(z) \log \frac{p(z)}{q(z)} dz$$

$$p(z|x) = \min_{p(z|x)} \{-I(X, Z) + \lambda KL(p(z)||q(z))\}$$

$$= \min_{p(z|x)} \left\{ - \iint p(z|x) \tilde{p}(x) \log \frac{p(z|x)}{p(z)} dx dz + \lambda \int p(z) \log \frac{p(z)}{q(z)} dz \right\}$$

$$= \min_{p(z|x)} \left\{ \iint p(z|x) \tilde{p}(x) \left[ -(1 + \lambda) \log \frac{p(z|x)}{p(z)} + \lambda \log \frac{p(z|x)}{q(z)} \right] dx dz \right\}$$

$$= \min_{p(z|x)} \left\{ -\beta \cdot I(X, Z) + \gamma \cdot \mathbb{E}_{x \sim \tilde{p}(x)} [KL(p(z|x)||q(z))] \right\}$$

- Mutual Information

$$\begin{aligned} I(X, Z) &= \iint p(z|x)\tilde{p}(x) \log \frac{p(z|x)\tilde{p}(x)}{p(z)\tilde{p}(x)} dx dz \\ &= KL(p(z|x)\tilde{p}(x) \| p(z)\tilde{p}(x)) \end{aligned}$$

- From f-GAN

$$\mathcal{D}_f(P \| Q) = \max_T \left( \mathbb{E}_{x \sim p(x)} [T(x)] - \mathbb{E}_{x \sim q(x)} [g(T(x))] \right)$$

$$JS(p(z|x)\tilde{p}(x), p(z)\tilde{p}(x))$$

$$= \max_T \left( \mathbb{E}_{x \sim p(z|x)\tilde{p}(x)} [\log \sigma(T(x, z))] + \mathbb{E}_{x \sim p(z)\tilde{p}(x)} [\log(1 - \sigma(T(x, z)))] \right)$$

- Deep InfoMax

$$p(z|x) = \min_{p(z|x)} \left\{ -\beta \cdot JS(p(z|x)\tilde{p}(x), p(z)\tilde{p}(x)) + \gamma \cdot \mathbb{E}_{x \sim \tilde{p}(x)} [KL(p(z|x) \| q(z))] \right\}$$

$$p(z|x), T(x, z)$$

$$= \min_{p(z|x), T(x, z)} \left\{ -\beta \cdot \left( \mathbb{E}_{x \sim p(z|x)\tilde{p}(x)} [\log \sigma(T(x, z))] + \mathbb{E}_{x \sim p(z)\tilde{p}(x)} [\log(1 - \sigma(T(x, z)))] \right) + \gamma \cdot \mathbb{E}_{x \sim \tilde{p}(x)} [KL(p(z|x) \| q(z))] \right\}$$

Discriminator

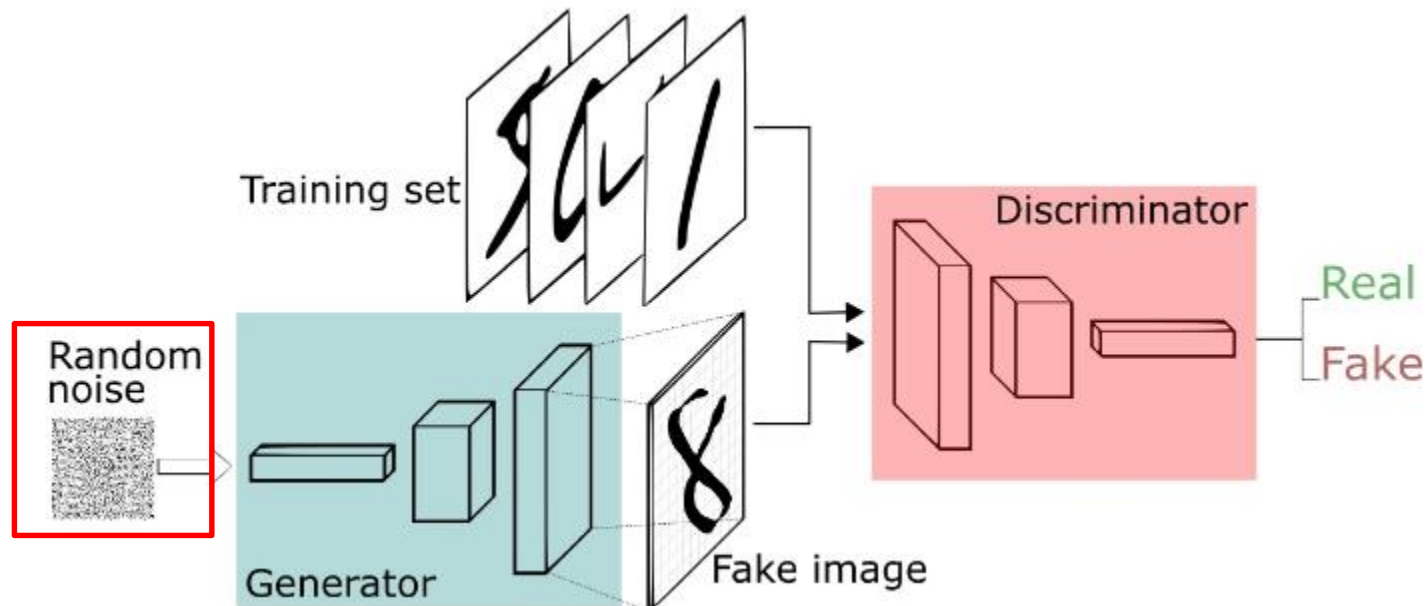
$$KL(p(z|x) \| q(z)) = \frac{1}{2} \sum_{k=1}^d \left( \mu_{(k)}^2(x) + \sigma_{(k)}^2(x) - \ln \sigma_{(k)}^2(x) - 1 \right)$$



# Disentangled Representation

- GAN

$$\min_G \max_D V(D, G) = E_{p(\mathbf{x})} [\log(D(\mathbf{x}))] + E_{p(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$



- InfoGAN

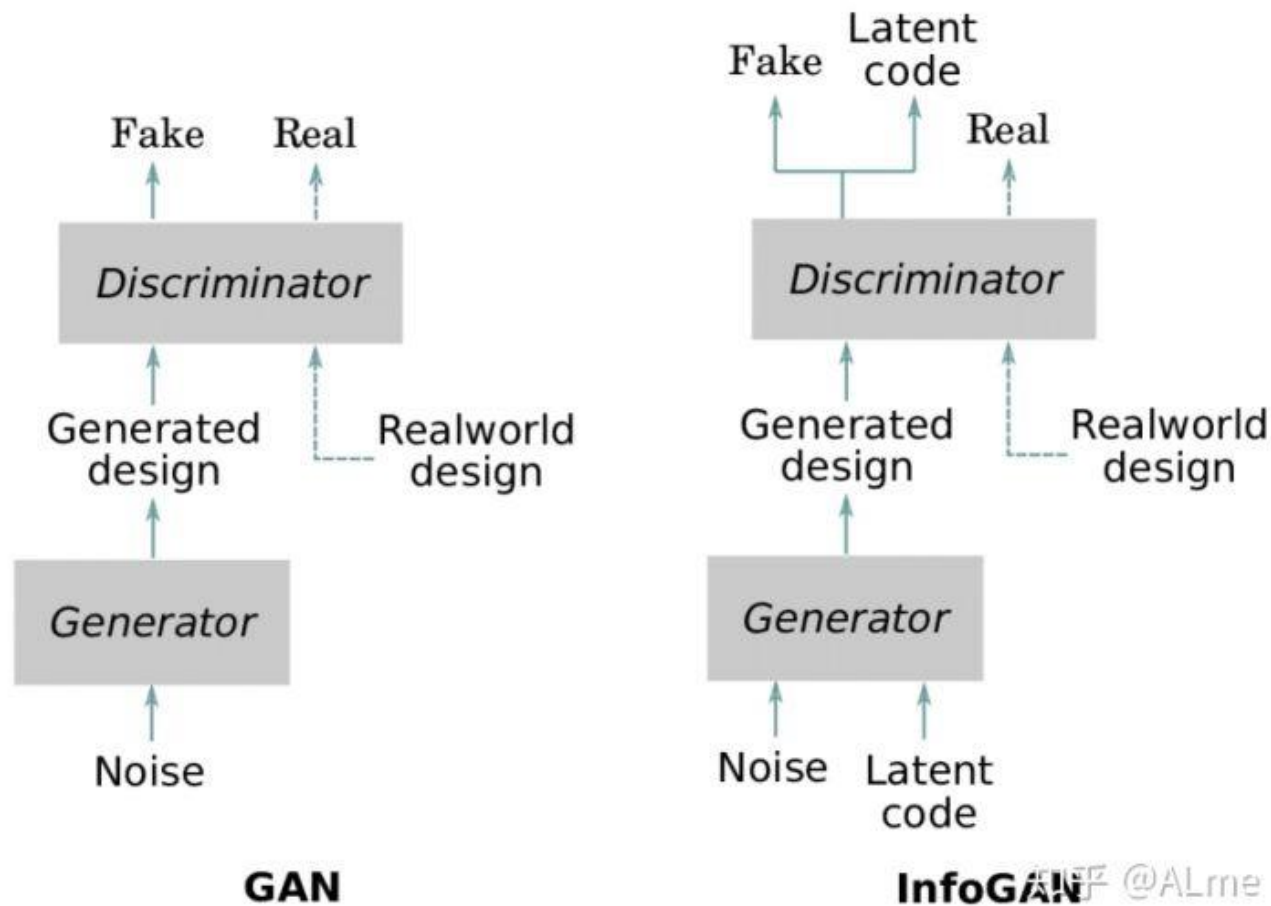
$$\min_G \max_D V_I(D, G) = V(D, G) - \lambda I(c; G(z, c))$$

$$\begin{aligned} I(c; G(z, c)) &= H(c) - H(c|G(z, c)) \\ &= \mathbb{E}_{x \sim G(z, c)} [\mathbb{E}_{c' \sim P(c|x)} [\log P(c'|x)]] + H(c) \\ &= \mathbb{E}_{x \sim G(z, c)} [\underbrace{D_{\text{KL}}(P(\cdot|x) \parallel Q(\cdot|x))}_{\geq 0} + \mathbb{E}_{c' \sim P(c|x)} [\log Q(c'|x)]] + H(c) \\ &\geq \mathbb{E}_{x \sim G(z, c)} [\mathbb{E}_{c' \sim P(c|x)} [\log Q(c'|x)]] + H(c) \end{aligned}$$

$$L_I(G, Q) = \mathbb{E}_{c \sim P(c), x \sim G(z, c)} [\log Q(c|x)] + H(c)$$

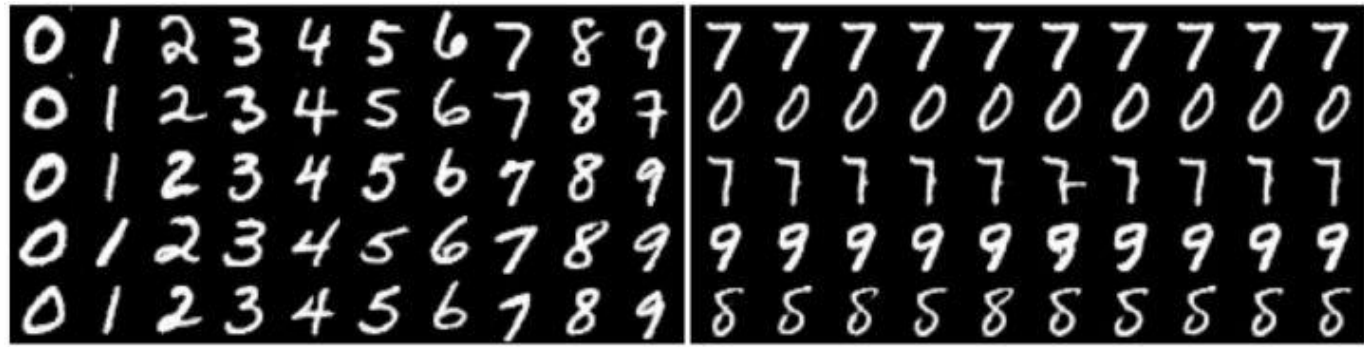
Discriminator

- Architecture



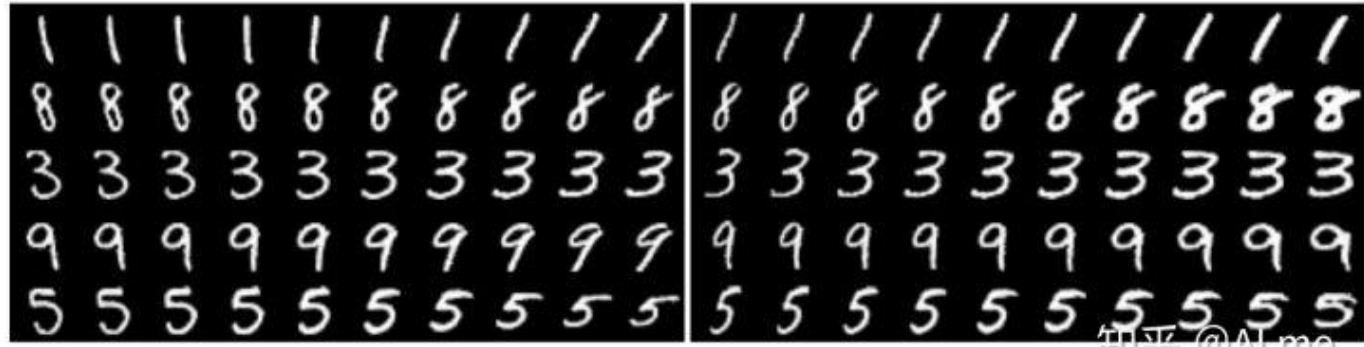


- Disentangled Representation



(a) Varying  $c_1$  on InfoGAN (Digit type)

(b) Varying  $c_1$  on regular GAN (No clear meaning)



(c) Varying  $c_2$  from  $-2$  to  $2$  on InfoGAN (Rotation)

(d) Varying  $c_3$  from  $-2$  to  $2$  on InfoGAN (Width)

知乎 @ALme

- VAE

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))$$

- BetaVAE

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))$$

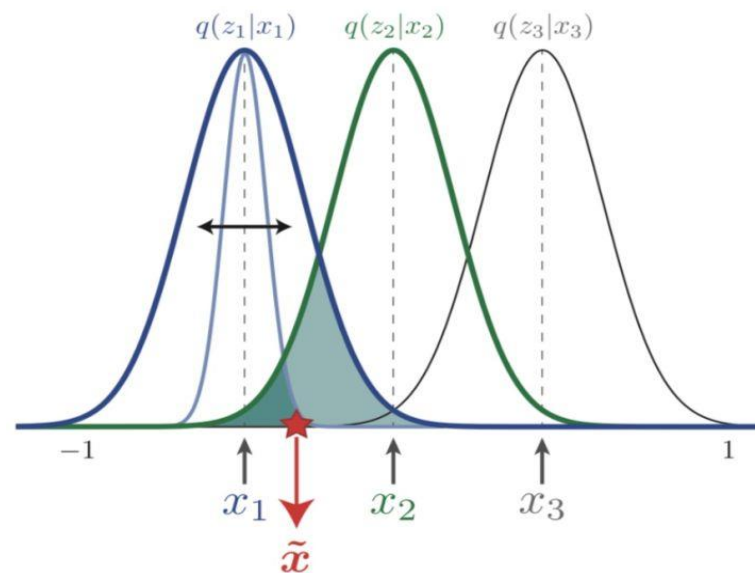
Standard Gaussian distribution

- Information Bottleneck

$$\max [I(Z; Y) - \beta I(X; Z)]$$

- Understanding BetaVAE

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))$$



# 4.3. Standard Gaussian Distribution

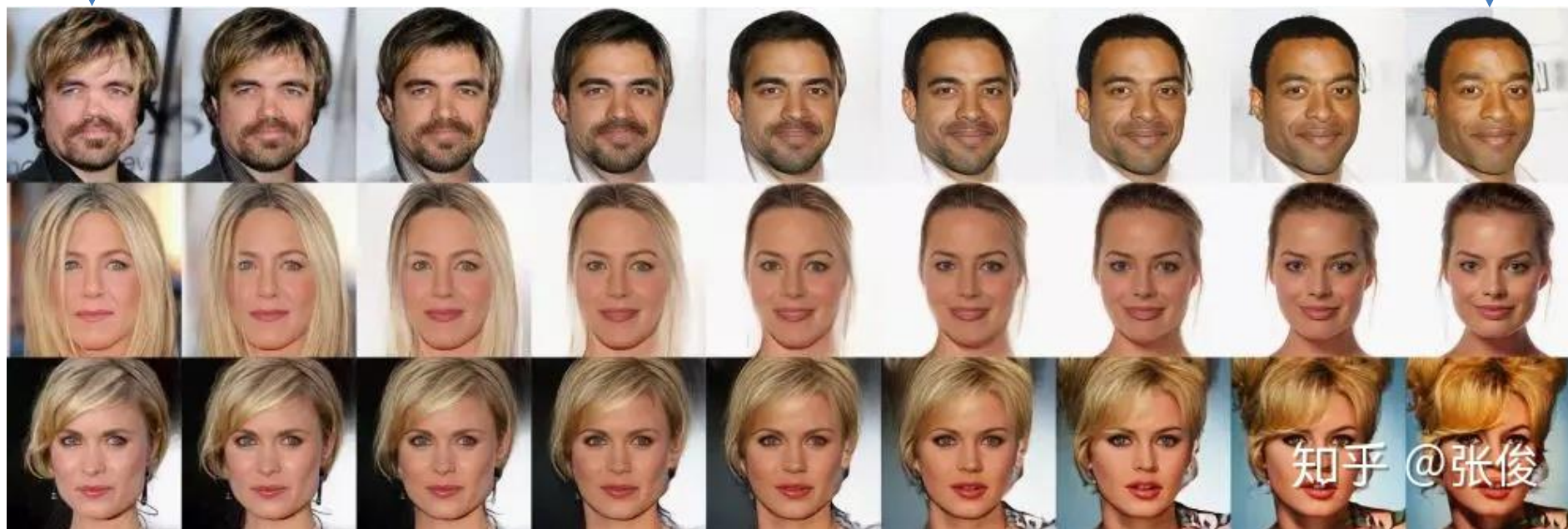


- Decomposition

$$P(x, y) = p(x)p(y)$$

- Interpolation

$$z_1 + z_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$





# Generalization Bound

- Generalization Error

$$E(h; \mathcal{D}) = P_{\mathbf{x} \sim \mathcal{D}}(h(\mathbf{x}) \neq y)$$

- Empirical error

$$\hat{E}(h; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(h(\mathbf{x}_i) \neq y_i)$$

- Generalization bound

$$E(h) \leq \epsilon$$

- Probably Approximately Correct (PAC)

$$P(E(h) \leq \epsilon) \geq 1 - \delta, \quad 0 < \epsilon, \delta < 1.$$

- Sample Complexity

$$m \geq \text{poly}(1/\epsilon, 1/\delta, \text{size}(\mathbf{x}), \text{size}(c))$$

- Properly PAC learnable

$$\mathcal{H} = \mathcal{C}$$

- Infinity hypothesis space

$$|\mathcal{H}| = \infty$$

- If  $|\mathcal{H}| \neq \infty$  and  $c \in \mathcal{H}$

$$\begin{aligned} P(h(\mathbf{x}) = y) &= 1 - P(h(\mathbf{x}) \neq y) \\ &= 1 - E(h) \\ &< 1 - \epsilon . \end{aligned}$$

$$\begin{aligned} P((h(\mathbf{x}_1) = y_1) \wedge \dots \wedge (h(\mathbf{x}_m) = y_m)) &= (1 - P(h(\mathbf{x}) \neq y))^m \\ &< (1 - \epsilon)^m . \end{aligned}$$

$$\begin{aligned} P(h \in \mathcal{H} : E(h) > \epsilon \wedge \hat{E}(h) = 0) &< |\mathcal{H}|(1 - \epsilon)^m \\ &< |\mathcal{H}|e^{-m\epsilon} \leq \delta , \end{aligned}$$

$$m \geq \frac{1}{\epsilon} \left( \ln |\mathcal{H}| + \ln \frac{1}{\delta} \right)$$



- If  $|\mathcal{H}| \neq \infty$  &  $c \notin \mathcal{H}$

$$P(E(h) - \min_{h' \in \mathcal{H}} E(h') \leq \epsilon) \geq 1 - \delta$$

- If  $|\mathcal{H}| = \infty$

$$\text{VC}(\mathcal{H}) = \max\{m : \Pi_{\mathcal{H}}(m) = 2^m\}$$

$$P\left(E(h) - \hat{E}(h) \leq \sqrt{\frac{8d \ln \frac{2em}{d} + 8 \ln \frac{4}{\delta}}{m}}\right) \geq 1 - \delta.$$

- Low MI  $\Rightarrow$  Generalization

$$|\mathbf{E}[\phi_T - \mu_T]| \leq \sigma \sqrt{2I(T; \phi)},$$

$$|\mathbf{E}[\phi_T] - \mathbf{E}[\mu_T]| \leq \sigma \sqrt{\frac{2I(T; \phi)}{n}}.$$

$$\mathbf{E}[L(\hat{f}) - \hat{L}(\hat{f})] \leq \sqrt{\frac{I(\hat{f}(\mathbf{x}); \mathbf{Y})}{2n}}.$$

- $I(T(X), X) \rightarrow \infty$ , but strong generalization holds
- CMI is finite

$$CMI_D(T) = I(T(\tilde{X}_S); S | \tilde{X})$$

$$CMI_D(T) \leq H(S) = \log 2^n = n$$

- Low CMI  $\Rightarrow$  Generalization
- CMI is bounded by VC dimension and Compression Schemes
- CMI is bounded from  $(\epsilon, \delta)$ -DP (via TV Stability)

Thank  
you

