

((Alireza Gargoori Motlagh - 98102176))

((Homework 0 - Reinforcement Learning))

Estimator's Variance

$$1 \quad \text{Bias}_{\theta}(\omega) = \mathbb{E}_{\theta}[\omega] - \theta \quad , \quad \text{Var}_{\theta}(\omega) = \mathbb{E}_{\theta}\left[(\omega - \mathbb{E}_{\theta}[\omega])^2\right]$$

$$\text{MSE}(\omega, \theta) = \mathbb{E}_{\theta}\left[(\omega - \theta)^2\right] = \mathbb{E}_{\theta}\left[\left(\omega - \mathbb{E}_{\theta}[\omega] + \mathbb{E}_{\theta}[\omega] - \theta\right)^2\right] =$$
$$\mathbb{E}_{\theta}\left[(\omega - \mathbb{E}_{\theta}[\omega])^2 + 2(\omega - \mathbb{E}_{\theta}[\omega])(\mathbb{E}_{\theta}[\omega] - \theta) + (\mathbb{E}_{\theta}[\omega] - \theta)^2\right] =$$

$$\boxed{\mathbb{E}_{\theta}\left[(\omega - \mathbb{E}_{\theta}[\omega])^2\right]} + 2 \boxed{\mathbb{E}_{\theta}\left[(\omega - \mathbb{E}_{\theta}[\omega])(\mathbb{E}_{\theta}[\omega] - \theta)\right]} + \boxed{\mathbb{E}_{\theta}\left[(\mathbb{E}_{\theta}[\omega] - \theta)^2\right]} =$$

Var _{θ} (ω) \star $\star\star$

$$\star \quad \mathbb{E}\left[(\omega - \mathbb{E}_{\theta}[\omega])(\mathbb{E}_{\theta}[\omega] - \theta) | \theta\right] = (\mathbb{E}_{\theta}[\omega] - \theta) \mathbb{E}_{\theta}[\omega - \mathbb{E}_{\theta}[\omega]] =$$
$$(\mathbb{E}_{\theta}[\omega] - \theta) \underbrace{(\mathbb{E}_{\theta}[\omega] - \mathbb{E}_{\theta}[\mathbb{E}_{\theta}[\omega]])}_{= \mathbb{E}_{\theta}[\omega]} = (\mathbb{E}_{\theta}[\omega] - \theta) \times 0 = 0$$

$$\star\star \quad \mathbb{E}_{\theta} \left[(\mathbb{E}_{\theta}[\omega] - \theta)^2 \right] = \mathbb{E} \left[(\mathbb{E}_{\theta}[\omega] - \theta)^2 | \theta \right] = (\mathbb{E}_{\theta}[\omega] - \theta)^2 = (\text{Bias}_{\theta}(\omega))^2$$

$\xrightarrow{\text{**}}$

$$\text{MSE}(\omega, \theta) = \mathbb{E}_{\theta}[(\omega - \theta)^2] = \text{Var}_{\theta}(\omega) + (\text{Bias}_{\theta}(\omega))^2$$

So, to have a low error (low MSE), we must have a low variance (good precision) & low bias (good accuracy), which is quite not possible because of bias-variance tradeoff.

2

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n f(x_i) , \quad X_i \stackrel{\text{iid}}{\sim} P(X) , \quad \theta = \mathbb{E}[f(X)]$$

$$\text{Bias}_{\theta}(\hat{\theta}_n) = \mathbb{E}_{\theta}[\hat{\theta}_n] - \theta = \mathbb{E}_{\theta} \left[\frac{1}{n} \sum_{i=1}^n f(x_i) \right] - \theta \stackrel{\text{linearity}}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\theta}[f(x_i)] - \theta =$$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\theta}[f(x_i)] - \theta \stackrel{\text{iid}}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\theta}[f(X)] - \theta =$$

$$\frac{1}{n} \times (n\theta) - \theta = \theta - \theta = 0 \Rightarrow \boxed{\hat{\theta}_n \text{ is an unbiased estimator for } \theta.}$$

$$3 \quad f_1(x) = 1 + \left(\frac{x}{2}\right)^2, \quad f_2(x) = \left(\frac{x}{4}\right)^{10}, \quad X \sim \text{Uniform}[0, 4]$$

$$\mathbb{E}_P[f_1(x)] = \int_{-\infty}^{\infty} f_1(x) f_X(x) dx = \frac{1}{4} \int_0^4 \left(1 + \left(\frac{x}{2}\right)^2\right) dx = \frac{7}{3} \approx 2.33$$

$$\mathbb{E}_P[f_2(x)] = \int_{-\infty}^{\infty} f_2(x) f_X(x) dx = \frac{1}{4} \int_0^4 \left(\frac{x}{4}\right)^{10} dx = \frac{1}{11} \approx 0.091$$

The code I wrote produced: $\hat{\theta}_1 = 2.144$, $\hat{\theta}_2 = 0.047$

$$\text{Bias}_{\hat{\theta}_1}(\hat{\theta}_1) = -0.19$$

$$\text{Bias}_{\hat{\theta}_2}(\hat{\theta}_2) = -0.044$$

$$\text{Var}_{\hat{\theta}_1}(\hat{\theta}_1) = 1.003$$

$$\text{Var}_{\hat{\theta}_2}(\hat{\theta}_2) = 0.017$$

As we can see, the bias and variance of this estimator is greater (in magnitude) for $f_1(x)$ than $f_2(x)$. The reason is that $f_1(x)$ is greater than $f_2(x)$ in magnitude and hence, its estimation has a larger bias & variance; since Monte Carlo estimator's variance is proportional to the magnitude of the estimating function.

$$\hat{\theta} = \mathbb{E}_q \left[\underbrace{\frac{P(\alpha)}{q(\alpha)} * \ell(\alpha)}_{\hat{\theta}} \right] = \sum_{\alpha} q(\alpha) \frac{P(\alpha)}{q(\alpha)} \ell(\alpha) = \sum_{\alpha} P(\alpha) \ell(\alpha) = \mathbb{E}_P [\ell(\alpha)]$$

$$\Rightarrow \text{Bias}(\hat{\theta}) = \mathbb{E}_q \left[\underbrace{\frac{P(\alpha)}{q(\alpha)} * \ell(\alpha)}_{\hat{\theta}} \right] - \underbrace{\mathbb{E}_P [\ell(\alpha)]}_{\theta} = 0 : \text{unbiasedness}$$

($q(\alpha) = 0$ must imply $P(\alpha) \ell(\alpha) = 0$)

$$\hat{\theta}_q = \frac{1}{n} \sum_{i=1}^n \frac{P(X_i) \ell(X_i)}{q(X_i)}, \quad \hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n \ell(X_i)$$

$$\text{Var}_q (\hat{\theta}_q) = \mathbb{E}_q \left[\left(\frac{P(X)}{q(X)} \ell(X) \right)^2 \right] - \left(\mathbb{E}_q \left[\frac{P(\alpha)}{q(\alpha)} * \ell(\alpha) \right] \right)^2 \Rightarrow$$

$$\text{Var}_q (\hat{\theta}_q) = \sum_{\alpha} q(\alpha) \left(\frac{P(\alpha)}{q(\alpha)} \ell(\alpha) \right)^2 - \theta^2$$

$$\text{Var}_P (\hat{\theta}_n) = \mathbb{E}_P [\ell^2(X)] - \left(\mathbb{E}_P [\ell(X)] \right)^2 = \sum_{\alpha} P(\alpha) \ell^2(\alpha) - \theta^2 \Rightarrow$$

$$\text{Var}_P (\hat{\theta}_n) = \sum_{\alpha} P(\alpha) \ell^2(\alpha) - \hat{\theta}^2$$

So, the reduction in the Variance is :

$$\text{Var}_P(\hat{\theta}_n) - \text{Var}_q(\hat{\theta}_q) = \sum_{\alpha} f^2(\alpha) P(\alpha) - \sum_{\alpha} f^2(\alpha) \frac{P^2(\alpha)}{q(\alpha)} = \sum_{\alpha} f^2(\alpha) \left(1 - \frac{P(\alpha)}{q(\alpha)}\right) P(\alpha)$$

\Rightarrow Reduced Var = $\sum_{\alpha} f^2(\alpha) \left(1 - \frac{P(\alpha)}{q(\alpha)}\right) P(\alpha)$

5. minimize $\sum_{\alpha} f^2(\alpha) \left(1 - \frac{P(\alpha)}{q(\alpha)}\right) P(\alpha)$
 s.t. $\sum_{\alpha} q(\alpha) = 1$

$\rightarrow L(q, \lambda) = \sum_{\alpha} f^2(\alpha) \left(1 - \frac{P(\alpha)}{q(\alpha)}\right) P(\alpha) + \lambda \left(\sum_{\alpha} q(\alpha) - 1\right)$

$$\frac{\partial L(q, \lambda)}{\partial q} = 0 \implies f^2(\alpha) \frac{P^2(\alpha)}{q^2(\alpha)} + \lambda = 0 \implies -\lambda q^2(\alpha) = f^2(\alpha) P^2(\alpha)$$

$$\rightarrow q^2(\alpha) = -\frac{1}{\lambda} f^2(\alpha) P^2(\alpha) \xrightarrow{q(\alpha) > 0} q^*(\alpha) = \frac{1}{|\lambda|} |f(\alpha)| P(\alpha) \rightarrow$$

$q^*(\alpha) \propto |f(\alpha)| P(\alpha)$

, $\lambda = \sum_{\alpha} |f(\alpha)| P(\alpha)$

We prove no other q has a lower variance than q^* :

$$q^*(x) = \frac{1}{\lambda} |f(x)| p(x), \quad \lambda = \sum_x |f(x)| p(x)$$

$$\text{Var}_{q^*} = \sum_x \frac{f(x)^2 p(x)}{q^*(x)} - \theta^2 = \lambda \sum_x \frac{f(x)^2 p(x)}{|f(x)| p(x)} - \theta^2$$

$$= \lambda \sum_x |f(x)| p(x) - \theta^2 = \left[\sum_x |f(x)| p(x) \right]^2 - \theta^2$$

$$= \left[\sum_x \frac{|f(x)| p(x)}{q(x)} q(x) \right]^2 - \theta^2 \stackrel{\text{Cauchy Schwarz}}{\leq} \sum_x \frac{f(x)^2 p(x)}{q(x)^2} q(x) - \theta^2$$

w.r.t. $q(x)$ (also $\sum_x q^2 \leq \sum_x q = 1$)

$$= \sum_x \frac{f(x)^2 p(x)}{q(x)} - \theta^2 = \text{Var}_q$$

6 $P = \mathbb{P}(X > 5), \quad X \sim \mathcal{N}(0, 1), \quad Q(x) = \begin{cases} e^{-(x-5)} & x \geq 5 \\ 0 & \text{o.w.} \end{cases}$

$$\mathbb{E}_P[f(x)] = \mathbb{E}_P[I_{X>5}] = \int_{-\infty}^{\infty} I_{X>5} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_5^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = 2.87 \times 10^{-7}$$

$$\mathbb{E}_Q \left[\frac{P(x)}{Q(x)} I_{x>5} \right] = \int_5^{\infty} e^{-(x-5)} \frac{\frac{1}{\sqrt{\pi}} e^{-\frac{x^2}{2}}}{e^{-(x-5)}} dx = \mathbb{E}_P \left[I_{x>5} \right] = 2.87 \times 10^{-7}$$

■ $\text{Var}_P \left(I_{x>5} \right) = \mathbb{E}_P \left[I_{x>5}^2 \right] - \left(\mathbb{E}_P \left[I_{x>5} \right] \right)^2 = \mathbb{E}_P \left[I_{x>5} \right] - \frac{\mathbb{E}_P \left[I_{x>5} \right]^2}{\mathbb{E}_P \left[I_{x>5} \right]} =$

$$2.87 \times 10^{-7} - 8.24 \times 10^{-14} \sim 10^{-7}$$

■ $\text{Var}_Q \left(\frac{P(x)}{Q(x)} I_{x>5} \right) = \mathbb{E}_Q \left[\left(\frac{P(x)}{Q(x)} I_{x>5} \right)^2 \right] - \left(\mathbb{E}_Q \left[\frac{P(x)}{Q(x)} I_{x>5} \right] \right)^2 =$

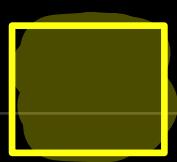
$$\int_{-\infty}^{\infty} Q(x) \frac{\frac{P(x)^2}{Q(x)} I^2_{x>5}}{Q(x)^2} dx - 8.24 \times 10^{-14} =$$

$$\int_5^{\infty} \frac{\frac{1}{\sqrt{\pi}} e^{-\frac{x^2}{2}}}{e^{-(x-5)}} dx - 8.24 \times 10^{-14} = 2.40 \times 10^{-13} - 8.24 \times 10^{-14} \sim 10^{-13}$$

As we can see, variance of the new estimator (importance sampling)

estimator) is about 10^{-13} , while the simple Monte Carlo estimator's

variance is about 10^{-7} ; So the variance is reduced significantly.



Markov Chain

$$1 \quad \mathbb{P}(X_n = s_n, X_{n-1} = s_{n-1}, \dots, X_0 = s_0) = \mathbb{P}(X_0 = s_0) \times \mathbb{P}(X_1 = s_1 | X_0 = s_0) \times$$

$$\mathbb{P}(X_2 = s_2 | X_1 = s_1, X_0 = s_0) \times \dots \times \mathbb{P}(X_n = s_n | X_{n-1} = s_{n-1}, X_{n-2} = s_{n-2}, \dots, X_0 = s_0)$$

markov property

$$= \mathbb{P}(X_0 = s_0) \times \mathbb{P}(X_1 = s_1 | X_0 = s_0) \times \mathbb{P}(X_2 = s_2 | X_1 = s_1) \times \dots \times \mathbb{P}(X_n = s_n | X_{n-1} = s_{n-1})$$

$$= \mathbb{P}(X_0 = s_0) \times \underbrace{P_x P_x \dots P}_n = P^{(0)} P^{(n)}$$

$$2 \quad \mathbb{P}(X_{t_1+t_2} | X_{F_1}, \dots, X_{F_n}) \stackrel{\text{Bayes}}{=} \frac{\mathbb{P}(X_{t_1+t_2}, X_{F_1}, X_{F_2}, \dots, X_{F_n})}{\mathbb{P}(X_{F_1}, X_{F_2}, \dots, X_{F_n})} =$$

$$\frac{P_x^{(0)} P^{(\min(F))} \times P^{(\max(F) - \min(F))} \times P^{(t_1+t_2) - \max(F)}}{P_x^{(0)} P^{(\min(F))} \times P^{(\max(F) - \min(F))}} = P^{(t_1+t_2) - \max(F)}$$

$$\Rightarrow \boxed{\mathbb{P}(X_{t_1+t_2} | X_{F_1}, \dots, X_{F_n}) = \mathbb{P}(X_{t_1+t_2} | X_{\max(F)})}$$

3 $P^{(n+m)} = \left[P_{ij}^{(n+m)} \right]_{ij}$ where $P_{ij}^{(n+m)} = P(X_{m+n} = j | X_0 = i)$

$$P_{ij}^{(n+m)} = P(X_{m+n} = j | X_0 = i) = \sum_k^* P(X_{m+n} = j, X_m = k | X_0 = i)$$

finite
n-states

$$= \sum_k \underbrace{P(X_{m+n} = j | X_m = k, X_0 = i)}_{\text{Chain rule}} P(X_m = k | X_0 = i)$$

$$= \sum_k \underbrace{P(X_{m+n} = j | X_m = k)}_{\text{Markov property}} P(X_m = k | X_0 = i)$$

$$= \sum_k P_{kj}^{(n)} P_{ik}^{(m)} = \sum_k P_{ik}^{(m)} P_{kj}^{(n)} \Rightarrow P^{(m+n)} = P^{(m)} P^{(n)}$$

4 $P^{(n)} = P^n = \begin{bmatrix} p_1 & 0 & \dots & 0 \\ 0 & p_2 & 0 & \dots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & p_n \end{bmatrix}^n = \begin{bmatrix} p_1^{(n)} & 0 & \dots & 0 \\ 0 & p_2^{(n)} & 0 & \dots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & p_n^{(n)} \end{bmatrix}$

$$P^{(o)} = \underbrace{\begin{bmatrix} \omega_1 & 0 & \dots & 0 \\ 0 & \omega_2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \omega_n \end{bmatrix}}_{\text{Weight matrix}} \begin{bmatrix} P_1^{(o)} \\ P_2^{(o)} \\ \vdots \\ P_n^{(o)} \end{bmatrix} = \begin{bmatrix} \omega_1 P_1^{(o)} \\ \omega_2 P_2^{(o)} \\ \vdots \\ \omega_n P_n^{(o)} \end{bmatrix}$$

$\omega := \text{weight matrix}$

$$\text{trace}(\omega) = \sum_{i=1}^n \omega_i = 1$$

$$\pi_{ss} = \lim_{n \rightarrow \infty} P^{(0)} P^{(n)} = \begin{bmatrix} w_1 P_1^{(0)} \\ w_2 P_2^{(0)} \\ \vdots \\ w_n P_n^{(0)} \end{bmatrix}^T \begin{bmatrix} P_1^{(n)} & 0 & \cdots & 0 \\ 0 & P_2^{(n)} & 0 & \cdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & P_n^{(n)} \end{bmatrix} = \begin{bmatrix} w_1 P_1^{(0)} P_1^{(n)} \\ w_2 P_2^{(0)} P_2^{(n)} \\ \vdots \\ w_n P_n^{(0)} P_n^{(n)} \end{bmatrix}$$

$$\Rightarrow \pi_{ss} = \begin{bmatrix} w_1 \pi_1 \\ w_2 \pi_2 \\ \vdots \\ w_n \pi_n \end{bmatrix}, \sum_{i=1}^n w_i = 1$$

So any convex combination of individual steady-states could be the steady-state for the combination of n independent markov chains ; hence , the π_{ss} distribution is not unique .

The long run behaviour of this markov chain would be determined by the coefficients of w , which indicate the importance (weight) of each individual markov chain . Therefore , $\pi_{ss} \sim \pi_{\arg \max(w)}$

5

$$\Pi_L(j) = \lim_{n \rightarrow \infty} P_{i,j}^{(n)} \quad \forall i, j$$

■ $\Pi_2(j) = \lim_{n \rightarrow \infty} P_{i,j}^{(n)} = \underset{3}{\lim_{n \rightarrow \infty}} \sum_k P_{i,k}^{(n-1)} P_{k,j} = \sum_k \lim_{n \rightarrow \infty} P_{i,k}^{(n-1)} P_{k,j}$

$$= \sum_k \Pi_1(k) P_{k,j} \quad \Rightarrow \quad \boxed{\Pi_L(j) = \sum_k \Pi_1(k) P_{k,j}, \quad \forall j}$$

which is equivalent to the stationary state : $\Pi = \Pi P$

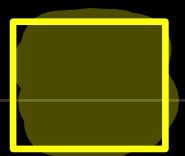
\Rightarrow The limiting distribution is always a stationary distribution.

■ $\Pi_1(j) = \sum_k \Pi_1(k) P_{k,j} \Rightarrow \Pi_2(j) = \sum_k \Pi_1(k) P_{k,j}^{(n)}$

$$\rightarrow \lim_{n \rightarrow \infty} \Pi_2(j) = \lim_{n \rightarrow \infty} \sum_k \Pi_1(k) P_{k,j}^{(n)} = \sum_k \underbrace{\lim_{n \rightarrow \infty} P_{k,j}^{(n)}}_{= \Pi_2(k)} \Pi_1(k)$$

$$= \sum_k \Pi_2(j) \Pi_1(k) = \Pi_2(j) \underbrace{\sum_k \Pi_1(k)}_{=1} = \Pi_2(j) \Rightarrow \boxed{\Pi_2^{(n)}(j) = \Pi_2(j)}$$

$\rightarrow \Pi^{(\sigma)}$ has no effect in limiting distribution on long run behaviour.



Information Theory

1 $H(X) = - \sum_x P(x) \log P(x) = - \mathbb{E}_X [\log P(X)]$

$$= \mathbb{E}_X [-\log P(X)] = \mathbb{E}_X \left[\log \frac{1}{P(X)} \right] \geq 0 \quad \text{since} \quad \log \frac{1}{P(X)} \geq 0$$

for $0 < P(X) \leq 1$

■ equality holds iff $\log \frac{1}{P(X)} = 0$ with probability 1 ; therefore X must be deterministic.

2 Since we want the maximum uncertainty (\equiv maximum entropy) , X must have a uniform distribution .

$X \sim \text{Uniform}(\frac{1}{M})$, where M is the number of outcomes of X RV :

$$\Rightarrow H(X) = - \sum_x P(x) \log P(x) = - \sum_x \frac{1}{M} \log \frac{1}{M} = -M \times \frac{1}{M} \times \log \frac{1}{M} = \log M$$

$$\Rightarrow H(X) = \log M$$

$$3 \quad H(X|Y) = - \sum_{x,y} P(x,y) \log P(x|y)$$

$$H(X,Y) = - \sum_{x,y} P(x,y) \log P(x,y) = - \sum_{x,y} P(x,y) \log (P(x|y)P(y)) =$$

$$- \sum_{x,y} P(x,y) \left(\log P(x|y) + \log P(y) \right) = - \sum_{x,y} P(x,y) \log P(x|y) -$$

$$\sum_{x,y} P(x,y) \log P(y) = H(X|Y) - \sum_y \log P(y) \underbrace{\sum_x P(x,y)}_{P(y)} =$$

$$H(X|Y) - \sum_y P(y) \log P(y) = H(X|Y) + H(Y)$$

$$H(X,Y) = - \sum_{x,y} P(x,y) \log (P(y|x) P(x)) = - \sum_{x,y} P(x,y) \log P(y|x) -$$

$$\sum_x \sum_y P(x,y) \log P(x) = H(Y|X) - \sum_x \log P(x) \underbrace{\sum_y P(x,y)}_{P(x)} =$$

$$H(Y|X) - \sum_x P(x) \log P(x) = H(Y|X) + H(X)$$

$$4 \quad D_{KL}(P || Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

Jensen inequality : $\mathbb{E}[g(x)] \geq g(\mathbb{E}[x])$ if $g(\cdot)$ is convex.

$\mathbb{E}[g(x)] \leq g(\mathbb{E}[x])$ if $g(\cdot)$ is concave.

$$-D_{KL}(P || Q) = -\sum_x P(x) \log \frac{P(x)}{Q(x)} = \sum_x P(x) \log \frac{Q(x)}{P(x)} \stackrel{\log(\cdot) \text{ is concave}}{\leq}$$

$$\log \sum_x P(x) \frac{Q(x)}{P(x)} = \log \sum_x Q(x) = -\log 1 = 0$$

$$\Rightarrow -D_{KL}(P || Q) \leq 0 \Rightarrow D_{KL}(P || Q) \geq 0$$

Jensen inequality is equality iff $g(\cdot)$ is affine or X is constant.

$$\Rightarrow 0 = D_{KL}(P || Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \geq 0 \stackrel{\begin{array}{l} \log \text{ is not} \\ \text{affine} \end{array}}{\Rightarrow} \frac{P(x)}{Q(x)} = \text{Const.} = 1$$

(since if $P(x) = CQ(x)$, due to the $\sum_x P(x) = 1 = C \sum_x Q(x) = C \Rightarrow C=1$)

$$\Rightarrow P(x) = Q(x)$$

$$5 \quad H(X) = - \sum_{\alpha} P(\alpha) \log P(\alpha) = -E_x [\log P(x)] = E_x [\log \frac{1}{P(x)}]$$

\log is a concave function ; Jensen inequality :

$$E_x [\log \frac{1}{P(x)}] \leq \log E_x \left[\frac{1}{P(x)} \right] = \log |\mathcal{X}| \quad \text{where } |\mathcal{X}| \text{ is the number of states of } X.$$

$$\Rightarrow H(X) \leq \log |\mathcal{X}|$$

Another way :

$$U(x) := \frac{1}{|\mathcal{X}|}, \quad \forall x \in \mathcal{X}$$

$$H(P) = - \sum_{\alpha} P(\alpha) \log P(\alpha) =$$

$$- \sum_{\alpha} P(\alpha) \log P(\alpha) + \sum_{\alpha} P(\alpha) \log U(\alpha) - \sum_{\alpha} P(\alpha) \log U(\alpha)$$

$$= - D_{KL}(P || U) - \sum_{\alpha} P(\alpha) \log \frac{1}{|\mathcal{X}|}$$

$$= - D_{KL}(P || U) + \log |\mathcal{X}| \sum_{\alpha} P(\alpha) = \log |\mathcal{X}| - D_{KL}(P || U)$$

$$\Rightarrow H(P) = \log |\mathcal{X}| - \underbrace{D_{KL}(P || U)}_{\geq 0} \leq \log |\mathcal{X}|$$

equality : $D_{KL}(P || U) = 0 \iff \boxed{P(\alpha) = U(\alpha) = \frac{1}{|\mathcal{X}|} : P \text{ is uniform}}$

6 We want to show that $H(X|Y) \leq H(X)$

$$I(X;Y) = D_{KL}(\hat{P}(X,Y) \parallel \hat{P}(X)\hat{P}(Y)) = \sum_{x,y} \hat{P}(x,y) \log \frac{\hat{P}(x,y)}{\hat{P}(x)\hat{P}(y)}$$

$$= \sum_{x,y} \hat{P}(x,y) \log \frac{\hat{P}(x|y)\hat{P}(y)}{\hat{P}(x)\hat{P}(y)} = \sum_{x,y} \hat{P}(x,y) \log \frac{\hat{P}(x|y)}{\hat{P}(x)} =$$

$$\sum_{x,y} \hat{P}(x,y) \log \hat{P}(x|y) - \sum_{x,y} \hat{P}(x,y) \log \hat{P}(x) = -H(X|Y) - \sum_x \log \hat{P}(x) \sum_y \hat{P}(x,y)$$

$$= -H(X|Y) - \sum_x \hat{P}(x) \log \hat{P}(x) = H(X) - H(X|Y)$$

Also, since $I(X;Y)$ is a KL divergence, it is ≥ 0 (as proved in 4):

$$I(X;Y) = H(X) - H(X|Y) \geq 0 \Rightarrow$$

$H(X) \geq H(X|Y)$

$$\text{7} \quad \blacksquare I(X, Y; Z) \geq I(X; Z) = H(X|Z) + H(Y|X, Z)$$

$$I(X, Y; Z) = \boxed{H(X, Y)} - \boxed{H(X, Y|Z)} = \\ = H(X) + H(Y|X)$$

$$\underbrace{[H(X) - H(X|Z)]}_{I(X; Z)} + \underbrace{[H(Y|X) - H(Y|X, Z)]}_{I(Z; Y|X)} \geq 0 \implies$$

$$I(X, Y; Z) - I(X; Z) \geq 0 \implies \boxed{I(X, Y; Z) \geq I(X; Z)}$$

equality: $I(Z; Y|X) = 0 \Rightarrow D_{KL}(\hat{P}(Z, Y|X) || P(Z|X)P(Y|X)) = 0$

$$\Rightarrow \boxed{\hat{P}(Z, Y|X) = \hat{P}(Z|X)\hat{P}(Y|X)}$$

$$\blacksquare H(X, Y|Z) \geq H(X|Z)$$

$$H(X, Y|Z) = H(X|Z) + \underbrace{H(Y|X, Z)}_{\geq 0} \geq H(X|Z)$$

equality: $H(Y|X, Z) = 0 \implies \boxed{Y = f(X, Z)}$

Bonus

8 ■ $I(X;Y) = 0, I(X;Y|Z) = 1$

$I(X;Y) = 0 \Rightarrow H(X) = H(X|Y) \Rightarrow X \& Y \text{ are independent.}$

$$I(X;Y,Z) = \cancel{I(X;Y)}_0 + I(X;Z|Y) = \underbrace{I(X;Z)}_{*=0} + \cancel{I(X;Y|Z)}_1$$

* For simplicity, we assume $I(X;Z)=0 \Rightarrow X \& Z \text{ are independent.}$

$$\Rightarrow X \sim \text{Bernoulli}(1/2), Z \sim \text{Bernoulli}(1/2), Y = \begin{cases} X & \text{if } Z=0 \\ 1-X & \text{if } Z=1 \end{cases}$$

| | | | | |
|---|-----|-----|-----|-----|
| | 00 | 01 | 10 | 11 |
| 0 | 1/4 | 0 | 0 | 1/4 |
| 1 | 0 | 1/4 | 1/4 | 0 |

(example taken from wikipedia)

■ $I(X;Y|Z)=0 \Rightarrow (X,Y) \text{ is independent of } Z.$

$$I(X;Y)=1 \Rightarrow \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} = 1$$

| | | | | |
|---|-----|-----|-----|-----|
| | 00 | 01 | 10 | 11 |
| 0 | 0 | 1/4 | 1/4 | 0 |
| 1 | 1/4 | 0 | 0 | 1/4 |



Probabilistic Models & Latent Variables

1 $p \sim$ true (unknown) distribution , $q_{\theta} \sim$ our estimation of p

$$D_{KL}(p||q) = \mathbb{E}_{x \sim p} \left[\log \frac{p(x)}{q_{\theta}(x)} \right] = \mathbb{E}_{x \sim p} [\log p(x)] - \mathbb{E}_{x \sim p} [\log q_{\theta}(x)]$$

$$\arg \min_{\theta} D_{KL}(p||q) = \underbrace{\arg \min_{\theta} \left(\mathbb{E}_{x \sim p} [\log p(x)] - \mathbb{E}_{x \sim p} [\log q_{\theta}(x)] \right)}_{\text{does not depend on } \theta}$$

$$= \arg \min_{\theta} - \mathbb{E}_{x \sim p} [\log q_{\theta}(x)] = \arg \max_{\theta} \mathbb{E}_{x \sim p} [\log q_{\theta}(x)]$$

Also, by using the Law of large numbers, we can say :

$$\arg \max_{\theta} \mathbb{E}_{x \sim p} [\log q_{\theta}(x)] = \arg \max_{\theta} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log q_{\theta}(x_i) =$$

$$\arg \max_{\theta} \log \prod_{i=1}^n q_{\theta}(x_i) = \arg \max_{\theta} \log q(\mathbf{x}|\theta) = \arg \max_{\theta} q(\mathbf{x}|\theta) = \hat{\theta}_{ML}$$

$$\Rightarrow \boxed{\arg \min_{\theta} D_{KL}(p||q_{\theta}) = \arg \max_{\theta} q(\mathbf{x}|\theta)}$$

$$2 \quad l(\theta; X) = \log \hat{P}_\theta(X) = \log \sum_z \hat{P}_\theta(X, z) = \log \sum_z \frac{\hat{P}_\theta(X, z)}{Q(z)} Q(z) =$$

$$\log \mathbb{E}_{Z \sim Q} \left[\frac{\hat{P}_\theta(X, z)}{Q(z)} \right] \stackrel{*}{\geqslant} \mathbb{E}_{Z \sim Q} \left[\log \frac{\hat{P}_\theta(X, z)}{Q(z)} \right]$$

* inequality holds since \log is concave (Jensen inequality)

$$\mathbb{E}_{Z \sim Q} \left[\log \frac{\hat{P}_\theta(X, z)}{Q(z)} \right] = \mathbb{E}_{Z \sim Q} \left[\log \hat{P}_\theta(X, z) \right] - \mathbb{E}_{Z \sim Q} \left[\log Q(z) \right] =$$

$$\mathbb{E}_{Z \sim Q} \left[\log \hat{P}_\theta(X, z) \right] + H(Q) = \mathcal{L}(\theta, Q)$$

$$\Rightarrow \boxed{l(\theta; X) = \log \hat{P}_\theta(X) \geqslant \mathcal{L}(\theta, Q) = \mathbb{E}_{Z \sim Q} \left[\log \hat{P}_\theta(X, z) \right] + H(Q)}$$

So, we have found a lower bound for log-likelihood parameterized by a distribution $Q(z)$ over Z (latent variables).

3 For a given θ , to minimize the distance between log-likelihood and lower bound, we must maximize $\mathcal{L}(\theta, Q)$:

$$\boxed{\begin{array}{l} \max \mathcal{L}(Q; \theta) \\ \text{s.t. } \sum_z Q(z) = 1 \end{array}}$$

(Constraint for a valid $Q(z)$)

$$\Rightarrow \text{Lagrange multipliers: } h(Q, \lambda) = \mathcal{L}(Q; \theta) + \lambda \left(\sum_z Q(z) - 1 \right)$$

$$= \mathbb{E}_{z \sim Q} [\log P_\theta(x, z)] - \mathbb{E}_{z \sim Q} [\log Q(z)] + \lambda \left(\sum_z Q(z) - 1 \right) =$$

$$\sum_z Q(z) \log P_\theta(x, z) - \sum_z Q(z) \log Q(z) + \lambda \left(\sum_z Q(z) - 1 \right)$$

$$\frac{\partial h(Q, \lambda)}{\partial Q(z)} = 0 \Rightarrow \sum_z \log P_\theta(x, z) - \left(\sum_z \log Q(z) + \sum_z \frac{Q(z)}{Q(z)} \right)$$

$$+ \lambda \sum_z 1 = 0 \xrightarrow[\text{out}]{} \log P_\theta(x, z) - \log Q(z) - 1 + \lambda = 0 \Rightarrow$$

$$\log Q(z) = \log P_\theta(x, z) + \lambda - 1 \implies$$

$$Q(z) = e^{\lambda - 1} P_\theta(x, z) \implies Q(z) \propto P_\theta(x, z) = P_\theta(z|x) \hat{P}(x)$$

$$Q(z) = e^{\lambda - 1} \underbrace{P_\theta(x)}_{\text{const.}} P_\theta(z|x) \implies Q^*(z) = P_\theta(z|x)$$

$$\Rightarrow \mathcal{L}(\theta, Q) = \mathbb{E}_{z \sim Q^*} \left[\log \frac{P_\theta(x, z)}{Q^*(z)} \right] = \sum_z Q^*(z) \log \frac{P_\theta(x, z)}{Q^*(z)}$$

$$= \sum_z \hat{P}_\theta(z|x) \log \frac{P_\theta(x, z)}{\hat{P}_\theta(z|x)} = \sum_z \hat{P}_\theta(z|x) \log \hat{P}_\theta(x) =$$

$$\log \hat{P}_\theta(x) \sum_z \hat{P}_\theta(z|x) = \log \hat{P}_\theta(x) = \mathcal{I}(\theta, x)$$

\implies The distance is 0 and the bound is tight.

4 As described, EM can be seen as E optimizing the lower bound of log-likelihood with respect to $Q(z)$ for a fixed θ which would be the posterior distribution $P(z|x)$, and M optimizing with respect to θ for a fixed Q of the previous step.

The problem with the above procedure would be the interactability

of the posterior $\underset{\theta}{P}(z|x)$; since :

$$\underset{\theta}{P}(z|x) = \frac{P_{\theta}(x|z) P_{\theta}(z)}{P_{\theta}(x)} = \frac{P_{\theta}(x|z) P_{\theta}(z)}{\sum_z P_{\theta}(x|z) P_{\theta}(z)}$$

The denominator of the above posterior is interactable; since it must be calculated for any possible z in the latent space!, which is almost impossible.

$$5 \quad \mathcal{L}(\theta, \phi) = \mathbb{E}_{z \sim Q_\phi} \left[\log \frac{\mathbb{P}_\theta(x, z)}{\mathbb{P}_\theta(z)} \right] + H(Q_\phi) = \mathbb{E}_{z \sim Q_\phi} \left[\log \frac{\mathbb{P}_\theta(x|z)}{Q_\phi(z)} \right]$$

• $\log \mathbb{P}_\theta(x) = \log \sum_z \mathbb{P}_\theta(x|z) \mathbb{P}_\theta(z) = \log \mathbb{E}_{z \sim \mathbb{P}_\theta} [\mathbb{P}_\theta(x|z)]$

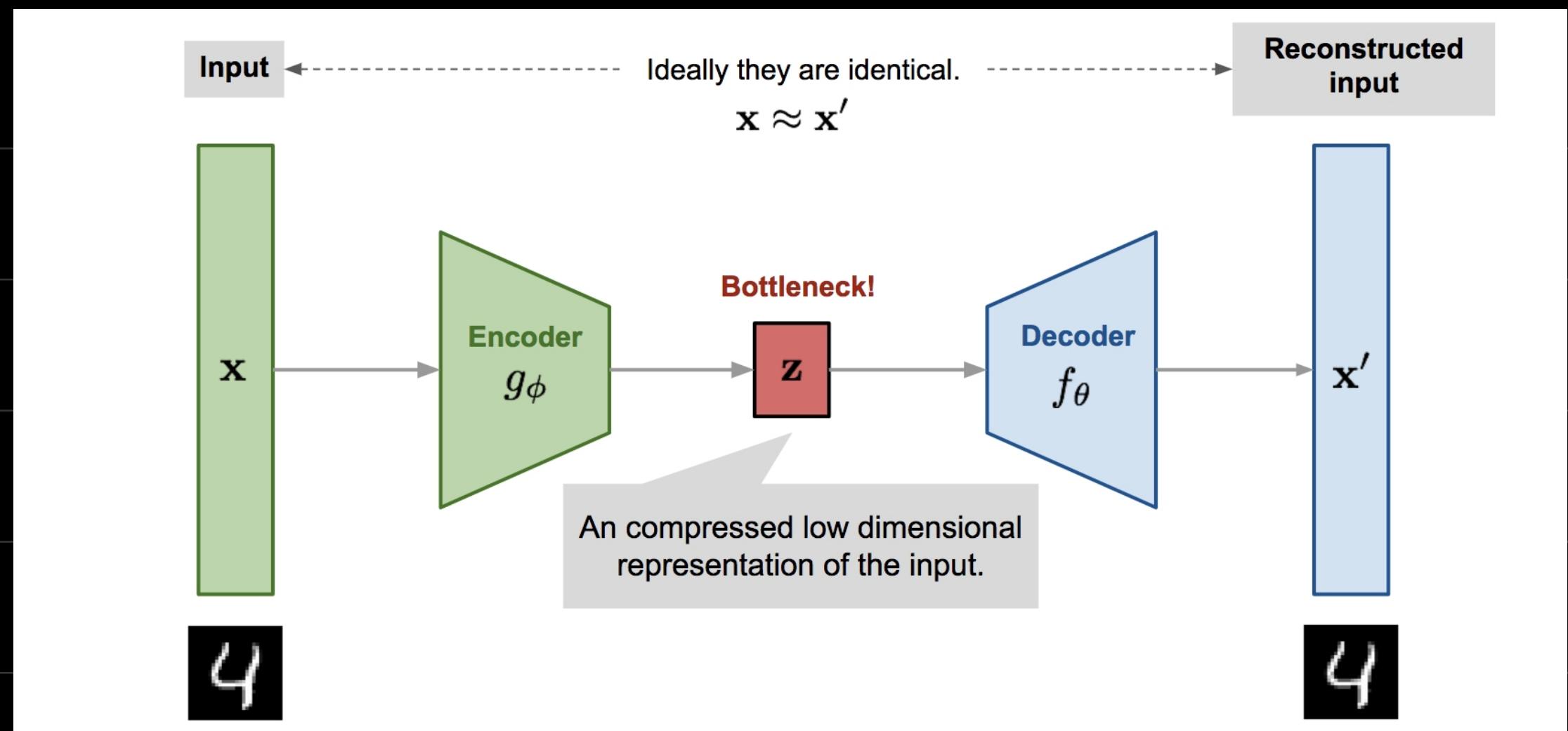
$$= \log \mathbb{E}_{z \sim \mathbb{P}_\theta} \left[\frac{\mathbb{P}_\theta(x|z)}{Q_\phi(z|x)} \right] = \log \mathbb{E}_{Q_\phi(z|x)} \left[\frac{\mathbb{P}_\theta(x|z) \mathbb{P}_\theta(z)}{Q_\phi(z|x)} \right]$$

$\stackrel{*}{\geqslant}$ Jensen $\mathbb{E} \left[\log \left(\frac{\mathbb{P}_\theta(x|z) \mathbb{P}_\theta(z)}{Q_\phi(z|x)} \right) \right] = \mathbb{E} \left[\log \left[\mathbb{P}_\theta(x|z) \right] \right] - \mathbb{E} \left[\log \left[Q_\phi(z|x) \right] \right]$

$$\mathbb{E} \left[\log \frac{Q_\phi(z|x)}{\mathbb{P}_\theta(z)} \right] = \mathbb{E} \left[\log \left(\mathbb{P}_\theta(x|z) \right) \right] - D_{KL}(Q_\phi(z|x) || \mathbb{P}_\theta(z))$$

$$= \mathcal{L}(\theta, \phi)$$

$\Rightarrow \boxed{\mathcal{L}(\theta, \phi) = \mathbb{E}_{Q_\phi(z|x)} [\log (\mathbb{P}_\theta(x|z))] - D_{KL}(Q_\phi(z|x) || \mathbb{P}_\theta(z))}$

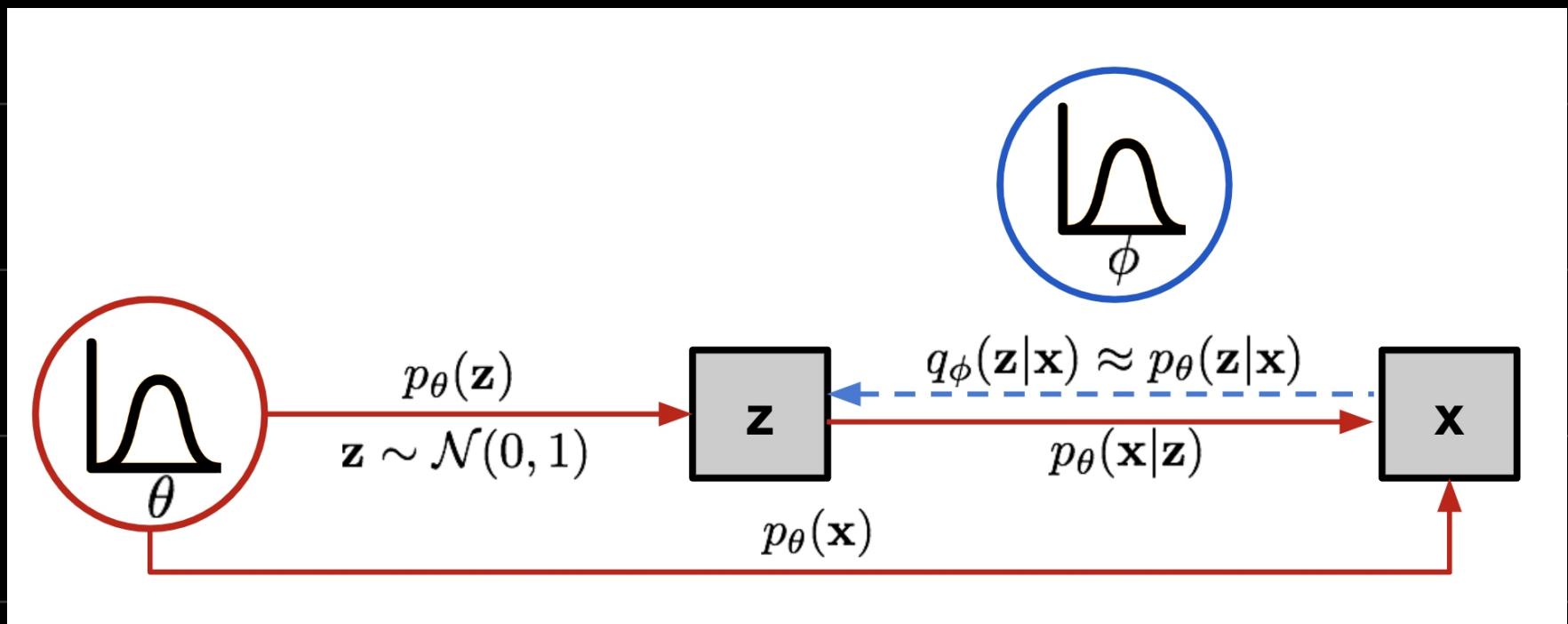


The training of the AE consists of minimizing a loss function through a deep neural network with encoder and decoder. A common loss could be MSE or ZINB (zero-inflated negative binomial, used for datasets with excessive dropout events, such as scRNA-seq data.)

$$\theta, \phi = \arg \min_{\theta, \phi} \mathcal{L}(f_\theta(g_\phi(\mathbf{x})), \mathbf{x}')$$

However, in the VAE, instead of mapping the input into a fixed vector in the latent space, we want to map it into a distribution.

As stated in the previous parts, maximizing the likelihood is equivalent to maximizing the ELBO (which is computable in contrast with intractable posterior); So, we estimate the underlying distribution with $q_\phi(z|x)$. The graphical model would be:



$$ELBO = \underbrace{E_{Q_\phi(z|x)} \left[\log \frac{P_\theta(x|z)}{p_\theta(z)} \right]}_{\text{likelihood}} - \underbrace{D_{KL} (Q_\phi(z|x) || P_\theta(z))}_{\text{Regularization : KL D}}$$

Minimizing the first term implies lower reconstruction error;

while the 2nd term is a regularizer, forcing the latent distribution to be as close

as possible to the $P_\theta(z)$, which is usually a standard normal distribution.

*The problem of not being able to backpropagate after sampling is addressed by reparam. trick)

In other words, the encoder takes the input x and outputs a mean μ and a variance σ^2 . Then, sampling occurs with $N(\mu, \sigma^2)$ through

$$\text{reparametrization trick: } z' \sim N(0, 1) \implies z = \beta \cdot z' + \mu \sim N(\mu, \sigma^2).$$

This random vector of latent space is passed through the decoder, which tries to reconstruct the original x . The KL divergence ensures the latent space distribution is close to the standard Normal distribution.

$$L(\theta, \phi) = \mathbb{E}_{z \sim Q_\phi(z|x)} \left[\log \frac{P_\theta(x|z)}{Q_\phi(z|x)} \right] - D_{KL}(Q_\phi(z|x) || P_\theta(z))$$

$$\mathbb{E}_{x \sim Q_\phi} \left[D_{KL}(Q_\phi(z|x) || P_\theta(z)) \right] = \sum_x Q_\phi(x) \sum_z Q_\phi(z|x) \log \frac{Q_\phi(z|x)}{P_\theta(z)} =$$

$$\sum_x Q_\phi(x) \sum_z Q_\phi(z|x) \log \frac{Q_\phi(z|x)}{P_\theta(z)} \times \frac{Q_\phi(z)}{Q_\phi(z)} = \sum_{x,z} Q_\phi(x, z) \log \frac{Q_\phi(x, z)}{Q_\phi(x) Q_\phi(z)}$$

$$\sum_{x,z} Q_\phi(x, z) \log \frac{Q_\phi(x)}{P_\theta(z)} = I(x; z) + D_{KL}(Q_\phi(z) || P_\theta(z))$$

$$\Rightarrow L(\theta, \phi) = \mathbb{E}_{Q_\phi(z|x)} \left[\log \frac{P_\theta(x|z)}{Q_\phi(z|x)} \right] - D_{KL}(Q_\phi(z) || P_\theta(z)) - I(x; z)$$

8

- By minimizing the KL divergence, we are implicitly minimizing the mutual information as well! This results in a poor reconstruction of the input, since the mutual information is an indicator of information of z in x and by reducing this quantity, we are encouraging our model to generate latent vectors independent of the input x . Hence, the reconstructed output would be more gaussian (more smooth) but has a higher error.

Solution :

$$\text{ELBO} := \mathbb{E}_{x \sim Q_\phi} \left[\mathbb{E}_{z \sim Q_\phi(z|x)} \left[\log P_\theta(x|z) \right] - \beta D_{\text{KL}}(Q_\phi(z|x) \parallel P_\theta(z)) \right]$$

, where β controls the mentioned tradeoff ; the lower the β , the higher the quality of reconstructed image & vice versa. (β -VAE)

- The noise variance estimation ϵ through maximizing the ELBO is biased.