(c) Powerful enough refers to if the neural network, parameterised by $\phi$, underlying $q_\phi$ can accurately represent the true posterior distribution. i.e. if $q_\phi(z|x) \approx p_\theta(z|x)$ sufficiently well ✓

$\exists \phi^* \text{ s.t. } \ldots \quad z_k \quad z_k \quad \text{for } k=1,\ldots,K$

In this case,

$$\underset{\theta,\phi}{\text{maximise}} \sum_{i=1}^{N} VLB_{\theta,\phi}^{(K)}(X_i)$$

close to equality by more powerful $q_\phi$ is
$$\widetilde{\approx} \underset{\theta}{\text{maximise}} \sum_{i=1}^{N} E_{z_1,\ldots,z_k \sim p_\theta(z|x)} \left[ \log \frac{1}{K} \sum_{k=1}^{K} \frac{p_\theta(x|z_k) p_z(z_k)}{p_\theta(z_k|x)} \right]$$ ✓

$$= \underset{\theta}{\text{maximise}} \sum_{i=1}^{N} \log\left( \frac{1}{K} \sum_{k=1}^{K} p_\theta(x) \right) = \underset{\theta}{\text{maximise}} \sum_{i=1}^{N} \log p_\theta(x) \quad \blacksquare ✓$$

2.(a) $\log p_\theta(X_i) = \log\left( E_{z \sim r_\lambda(z)}\left[ p_\theta(X_i|z) \right] \right)$

$$= \log\left( E_{z \sim q_\phi(z|X_i)}\left[ \frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right] \right)$$

$$\geq E_{z \sim q_\phi(z|X_i)}\left[ \log\left( \frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right) \right] \quad \text{by Jensen's inequality}$$

$$= VLB_{\theta,\phi,\lambda}(X_i) \quad \blacksquare ✓$$

(b) $\underline{\nabla} VLB(X_i) = \left( \nabla_\theta VLB(X_i), \nabla_\phi VLB(X_i), \nabla_\lambda VLB(X_i) \right)$

$$\nabla_\theta VLB(X_i) = \nabla_\theta \int \log\left( \frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right) q_\phi(z|X_i) dz$$

$$= \int \nabla_\theta\left( p_\theta(X_i|z) \right) \frac{1}{p_\theta(X_i|z)} q_\phi(z|X_i) dz + \underline{0}$$

$$= E_{z \sim q_\phi(z|X_i)}\left[ \nabla_\theta\left( \log\left( p_\theta(X_i|z) \right) \right) \right] ✓$$

$$\nabla_\phi VLB(X_i) = E_{z \sim q_\phi(z|X_i)}\left[ \left( \nabla_\phi \log q_\phi(z) \right) \log\left( \frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right) \right] \quad \text{by log-derivative trick for VAEs (Hw10.1)}$$

(with annotations $z|X_i$ and $\lambda$)

$$\nabla_\lambda VLB(X_i) = E_{z \sim q_\phi(z|X_i)}\left[ \nabla_\lambda\left( \log\left( r_\lambda(z) \right) \right) \right] \quad \text{by same logic as } \nabla_\theta ✓$$

So stochastic gradient of $VLB_{\theta,\phi,\lambda}(X_i)$ can be computed by:

$$\nabla_{\underline{\theta},\underline{\phi},\underline{\lambda}} VLB_{\theta,\phi,\lambda}(X_i) \approx \frac{1}{K} \sum_{k=1}^{K} \left( \nabla_\theta\left( \log\left( p_\theta(X_i|z_k) \right) \right), \left( \nabla_\phi \log(z_k) \right) \log\left( \frac{p_\theta(X_i|z_k) r_\lambda(z_k)}{q_\phi(z_k|X_i)} \right), \right.$$

$z_{i,k}|X_i$

number of elements in a batch of SG calculation

$$\left. \nabla_\lambda\left( \log\left( r_\lambda(z_k) \right) \right) \right) \quad \text{where } z_k \sim q_\phi(z|X_i)$$

$z_{i,k}$ (sample for each $X_i$ anew)

**2.(c)** $\nabla_\theta$ and $\nabla_\lambda$ same as above ✓ (Expectation distribution $q_\phi$ doesn't depend on $\theta$ or $\lambda$ so no need for the trick)

$\nabla_\phi$ evaluation changes if we use reparametrization vs. log-derivative before:

$$\nabla_\phi \, VLB(X_i) = \nabla_\phi \, \mathbb{E}_{z \sim q_\phi(z|X_i)} \left[ \log\left( \frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right) \right]$$

**Redo**  **2.(c)**
$$VLB_{\theta,\phi,\lambda}(X_i) = \underbrace{\mathbb{E}_{z \sim q_\phi(z|X_i)}[\log p_\theta(X_i|z)]}_{①} - \underbrace{\overbrace{D_{KL}\left( q_\phi(z|X_i) \| r_\lambda(z) \right)}^{+\mathbb{E}\left[\log\left(\frac{r_\lambda}{q_\phi}\right)\right]}}_{②} \qquad \text{(slide 89)}$$

- $① = \mathbb{E}_{z \sim q_\phi(z|X_i)} \left[ \log\left( (2\pi\sigma^2)^{-d/2} \exp\left( -\frac{1}{2}(X_i - f_\theta(z))^T \left(\frac{1}{\sigma^2}I\right)(X_i - f_\theta(z)) \right) \right) \right]$

   $d = \dim(f_\theta)$

   $= -\frac{1}{2\sigma^2} \mathbb{E}_{z \sim \mathcal{N}(\mu_\phi(X_i), \Sigma_\phi(X_i))} \left[ \| X_i - f_\theta(z) \|^2 \right] - \frac{d}{2}\log(2\pi\sigma^2)$

by reparametrization trick:

$$= -\frac{1}{2\sigma^2} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0,I)} \left[ \| X_i - f_\theta\left( \mu_\phi(X_i) + \Sigma_\phi^{1/2}(X_i)\varepsilon \right) \|^2 \right] - \frac{d}{2}\log(2\pi\sigma^2)$$

$\longrightarrow$ Apply $\nabla_\phi$: $-\frac{1}{2\sigma^2} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0,I)} \left[ \nabla_\phi \| \dots \|^2 \right]$   which can be evaluated + back-propagated

- $② = -\frac{1}{2}\left( tr\left( [diag(\lambda_2)]^{-1} \Sigma_\phi(X_i) \right) + (\lambda_1 - \mu_\phi(X_i))^T [diag(\lambda_2)]^{-1}(\lambda_1 - \mu_\phi(X_i)) - k + \log\left( \frac{\det(diag(\lambda_2))}{\det(\Sigma_\phi(X_i))} \right) \right)$

   by HW9.5

   $\longrightarrow$ Gradient $\nabla_\phi$ can then be calculated + back propagated