# Optimization of the loss function

As discussed, $\Theta^*, \phi^* = \underset{\Theta, \phi}{\text{argmin}} \, L(\Theta, \phi)$

In variational Baysian method, this loss function is known as the variational lower bound or evidence lower bound (ELBO). This 'lower bound' part comes from the fact that KL divergence is always non-negative & thus $L(\Theta, \phi)$ is the lower bound of $\log P_\Theta(x)$.

Recall. (ELBO proof).

$$\log P_\Theta(x) - D_{KL}\left[Q_\phi(z/x) \| P_\Theta(z/x)\right] = -L(\Theta, \phi)$$

And we know $D_{KL}\left[Q_{\phi}(z|x) \,||\, P_{\theta}(z|x)\right] \geq 0$ }

As a result

$$L(\theta, \phi) \leq \log P_{\theta}(x)$$ }

Therefore minimizing loss, we are maximizing the lower bound of the probability of generating real data samples.
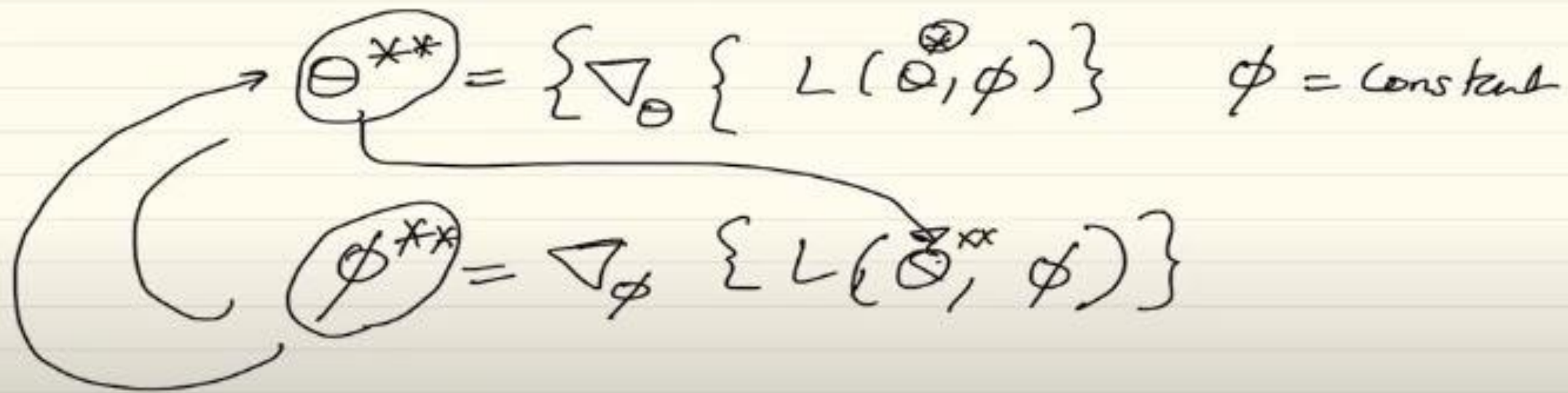
## REPARAMETRIZATION TRICK

Recall:                                    Needed during ( Back Propagation)

$$L(\theta, \phi) = - E_{z \sim Q_{\phi}(z|x)}\left[\log\left(P_{\theta}(x|z)\right)\right] +$$

$$\frac{1}{2}\sum_{K}\left\{\exp\left(\Sigma_{\phi}(x)\right) + \mu_{\phi}^{2}(x) - 1 - \Sigma_{\phi}(x)\right]$$

$$\Theta^*, \phi^* = \underset{\Theta, \phi}{argmin} \ L(\Theta, \phi)$$

Alternate optimization principle

$$\Theta^{**} = \{ \nabla_\Theta \{ L(\overset{\oslash}{\Theta}, \phi) \} \quad \phi = Constant$$

$$\phi^{**} = \nabla_\phi \{ L(\overset{\oslash**}{\Theta}, \phi) \}$$

Optimization is carried out with respect to both $\theta$ & $\phi$ to learn $Q_\phi(z/x)$ & $P_\theta(x/z)$ at the same time i.e.

$$\min_{\theta,\phi} L(\theta,\phi) = \min_{\theta,\phi} \left\{ -E_{z\sim Q_\phi(z/x)}\left[\log\left(P_\theta(x/z)\right)\right] + \frac{1}{2}\sum_K\left[\exp\left(\xi_\phi(x)\right) + \mu_\phi^2(x) - 1 - \xi_\phi(x)\right] \right\}$$

Alternate $i = 1,2,\ldots n$

(A) $\hat{\theta}_i = \nabla_\theta L(\theta,\phi)$

$$\hat{\theta}_i = \nabla_\theta\left\{ -E_{z\sim Q_\phi(z/x)}\left[\log P_\theta(x|z)\right] + \frac{1}{2}\sum_K\left[\exp\left(\xi_\phi(x)\right) + \mu_\phi^2(x) - 1 - \xi_\phi(x)\right] \right\}$$

$$\hat{\Theta}_i = \nabla_\Theta \mathcal{L}(\Theta, \phi)$$

$$\cong \frac{1}{L} \sum_{\ell=1}^{L} \underline{\nabla_\Theta \log P_\Theta(x/z^{(\ell)})} \quad \left\{ \begin{array}{c} \text{Monte Carlo} \\ \text{estimate} \end{array} \right\}.$$

where $z^{(\ell)} \sim Q_\phi(z/x)$

$$\hat{\Theta}_i = \nabla_\theta \mathcal{L}(\theta, \phi)$$

$$\hat{\Theta}_i \cong \frac{1}{L} \sum_{\ell=1}^{L} \nabla_\theta \log P_\theta(x | z^{(\ell)}) \qquad \left\{ \begin{array}{c} \text{Monte Carlo} \\ \text{estimate} \end{array} \right\}$$

$$\text{where} \quad z^{(\ell)} \sim Q_\phi(z|x)$$

$$\hat{\phi}_i = \nabla_\phi \left\{ \mathcal{L}(\theta, \phi) \right\}$$

$$(B) \quad \hat{\phi}_i = \nabla_\phi \left\{ L(\theta, \phi) \right\}$$

$$= \nabla_\phi \left\{ -E_{z \sim Q_\phi(z/x)} \left[ \log P_\theta(x|z) \right] + \right.$$

$$\left. \frac{1}{2} \sum_K \left[ \exp\left( \mathcal{E}_\phi(x) \right) + p_\phi^2(x) - 1 - \mathcal{E}_\phi(x) \right] \right\}$$

problem

This derivative $\nabla_\phi$ is harder to estimate because $\phi$ appears in the distribution with respect to which expectation is taken. i.e $\nabla_\phi E_{Q_\phi(y/x)}[f(z)] \neq E_{Q_\phi(y/x)}[\nabla_\phi f(z)]$

If we can somehow rewrite this expectation in such a way the $\phi$ appears inside the expectation then we can push the gradient inside the expectation i.e if we can write

$$E_{Q_\phi(z/x)}\left[f(z)\right] = E_{P(\epsilon)}\left[f\left(g_\phi(\epsilon,x)\right)\right]$$

such that $z = \boxed{g_\phi(\epsilon,x)}$ → Any linear transformation with $\epsilon \sim N(0,1)$
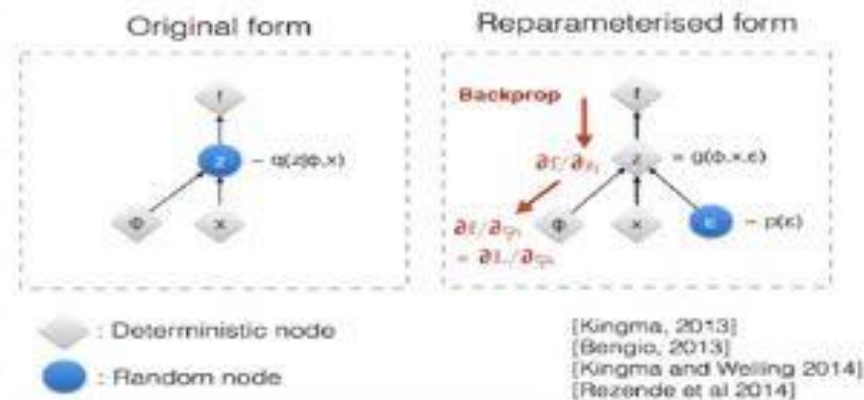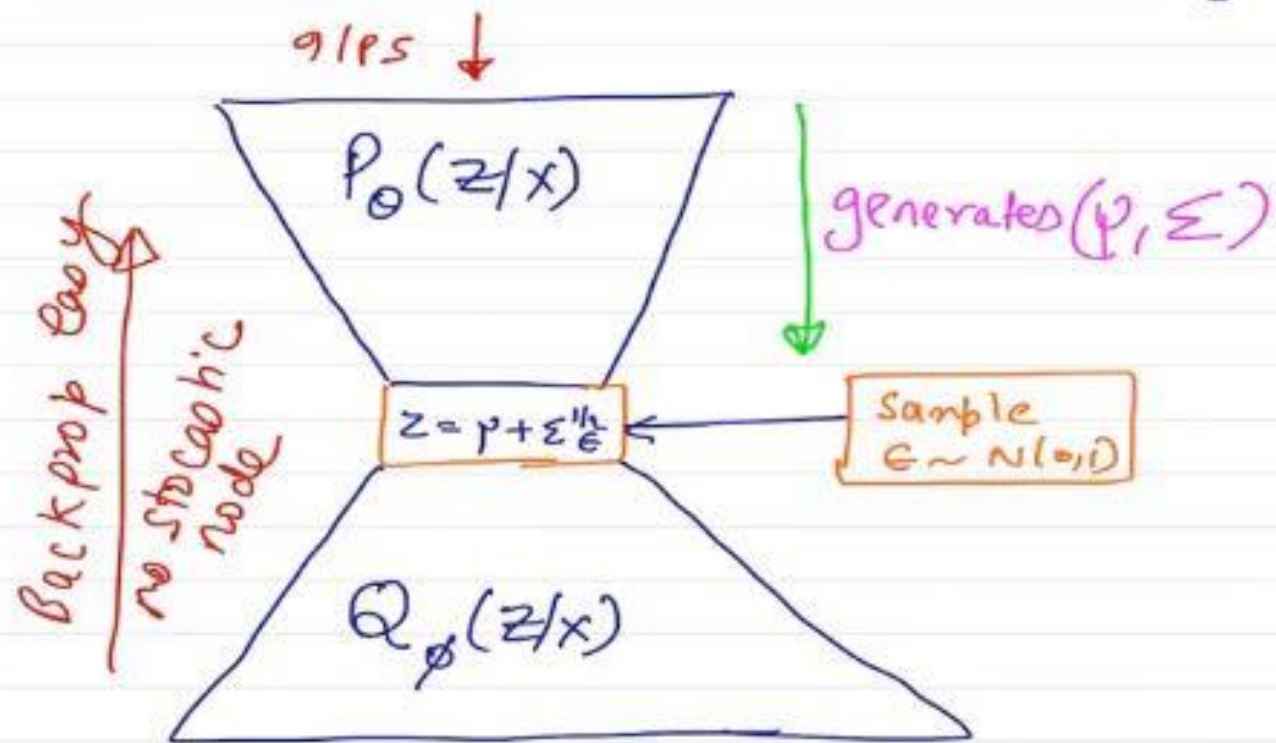
In our case $g_\phi(\epsilon,x) = P_\phi(x) + \epsilon \odot \Sigma_\phi^{1/2}(x) = z \sim N(P(x), \Sigma(x))$

Here, $N(P(x),\Sigma(x))$ is obtained from $N(0,1)$ using above linear transformation.

Instead of sampling $z \sim Q_\phi(z/x)$, we sample from $N(0,1)$

i.e $\quad \epsilon \sim N(0,1) \Rightarrow \epsilon \sim P(\epsilon)$

& then linear transform using $z = \mu_\phi(x) + \epsilon \odot \Sigma_\phi^{1/2}(x)$

to realise $N(\mu(x), \Sigma(x))$ defined earlier

steps $\downarrow$

$P_\theta(z/x)$

generates $(\mu, \Sigma)$

$z = \mu + \Sigma^{1/2}\epsilon$

Sample $\epsilon \sim N(0,1)$

$Q_\phi(z/x)$

Backprop Easy

no stochastic node



Original form

Reparameterised form

Backprop

$\partial L / \partial_{z_i}$   $z$   $= g(\phi, x, \epsilon)$

$\partial L / \partial_{z_i}$   $\phi$   $x$   $\epsilon$   $\sim p(\epsilon)$

$= \partial L / \partial_{z_i}$

◇ : Deterministic node

● : Random node

[Kingma, 2013]
[Bengio, 2013]
[Kingma and Welling 2014]
[Rezende et al 2014]

$$\hat{\phi}_i = \nabla_\phi \{ L(\theta, \phi) \}$$

$$= \nabla_\phi \left\{ -E_{z \sim Q_\phi(z/x)} \left[ \log P_\theta(x|z) \right] + \frac{1}{2} \sum_K \left[ exp\left( \Sigma_\phi(x) \right) + p_\phi^2(x) - 1 - \Sigma_\phi(x) \right] \right\}$$

problem

modified

$$\hat{\phi}_i = \nabla_\phi \{ L(\theta, \phi) \}$$

$$= \nabla_\phi \left\{ -E_{z^{(l)} \sim P(\epsilon)} \left[ \log P_\theta(x|z^{(l)}) \right] + \frac{1}{2} \sum_K \left[ exp\left( \Sigma_\phi(x) \right) + p_\phi^2(x) - 1 - \Sigma_\phi(x) \right] \right\}$$

$$\hat{\phi}_i = -\mathop{E}_{z^e \sim P(\epsilon)} \left[ \nabla_\phi \left( \log P_\theta \left( x | z^e \right) \right) \right] +$$

$$\nabla_\phi \left[ \frac{1}{2} \left( \sum_K \left[ \exp\left( \Sigma_\phi(x) \right) + \mu_\phi^2(x) - 1 - \Sigma_\phi(x) \right] \right) \right]$$

$$= \frac{-1}{S} \sum_{s=1}^{S} \cdot \left[ \log P_\theta \left( x | z^{(e)} \right) \right] +$$

Monte-Carlo estimate of Expectation.

where $z^{(e)} = \mu_\phi(x) + \epsilon \odot \sigma_\phi(x)$ & $\epsilon^{(e)} \sim N(0,1)$