1. Given the prior $p(z) \sim N(0, I)$ and the posterior approximation $q(z|x; \theta) \sim N(\mu_\theta(x), \sum_\theta(x))$, prove that $KL(q(z|x; \theta)||p(z))$ is tractable; that is, it can be the functions of $\mu_\theta(x)$ and $\sum_\theta(x)$, expressed as a closed-form expression. Both dimensions of multivariate Gaussian are $n$ where mean $\mu_\theta(x)$ and covariance matrix $\sum_\theta(x) = diag(\sigma_1^2, \ldots, \sigma_n^2)$ are functions of $x$ and the parameters $\theta$ of a neural network.

$$q(z|x;\theta) = \frac{1}{\sqrt{2\pi^n|\Sigma|}} \exp\left(-\frac{1}{2}(z-\mu_\theta(x))^\top \Sigma^{-1}(x)(z-\mu_\theta(x))\right) = N\left(\mu_\theta(x), \Sigma_\theta(x)\right)$$

$$p(z) = \frac{1}{\sqrt{2\pi^n|\Sigma|}} \exp\left(-\frac{1}{2}(z-0)^\top I^{-1}(z-0)\right) = N(0, I)$$

$$KL\left(q(z|x;\theta)||p(z)\right) = \int q(z|x;\theta) \log \frac{q(z|x;\theta)}{p(z)} dz$$

$$= \int q(z|x;\theta)\left[\log q(z|x;\theta) - \log p(z)\right] dz$$

$$= \int q(z|x;\theta)\left[-\frac{n}{2}\log(2\pi) - \frac{1}{2}\log(|\Sigma|) - \frac{1}{2\Sigma}(z-\mu)^\top(z-\mu) + \frac{n}{2}\log(2\pi) - \frac{1}{2}\log(|I|) + \frac{1}{2I}(z)^\top(z)\right]dz$$

$$= \int q(z|x;\theta)\left[\frac{1}{2}\log\frac{|I|}{|\Sigma|} + \frac{1}{2}\left(\frac{z^\top z}{I} - \frac{(z-\mu)^\top(z-\mu)}{\Sigma}\right)\right]dz$$

$$= E_q\left[\frac{1}{2}\log\frac{|I|}{|\Sigma|} + \frac{1}{2}\left(\frac{z^\top z}{I} - \frac{(z-\mu)^\top(z-\mu)}{\Sigma}\right)\right]$$

$$= \frac{1}{2}\log\frac{|I|}{|\Sigma|} + \frac{1}{2I}E_q[z^\top z] - \frac{1}{2\Sigma}E_q[(z-\mu)^\top(z-\mu)]$$

$$= \frac{1}{2}\log\frac{|I|}{|\Sigma|} + \frac{1}{2I}E_q[z^\top z] - \frac{1}{2\Sigma}\Sigma$$

$$= \frac{1}{2}\log\frac{|I|}{|\Sigma|} + \frac{1}{2I}E_q\left[(z-\mu+\mu+0)^\top(z-\mu+\mu+0)\right] - \frac{1}{2}tr(I_n)$$

$$= \frac{1}{2}\log\frac{|I|}{|\Sigma|} + \frac{1}{2I}\left[\underbrace{E_q[(z-\mu)^\top(z-\mu)]}_{} + 2\mu\underbrace{E_q[z-\mu]}_{0} + \mu^\top\mu\right] - \frac{1}{2}tr(I_n)$$

$$= \frac{1}{2}\log\frac{|I|}{|\Sigma|} + \frac{1}{2I}\left[\Sigma + \mu^\top\mu\right] - \frac{1}{2}n$$

$$= \frac{1}{2}\left[\log\frac{|I|}{|\Sigma|} + tr(I^{-1}\Sigma) + \mu^\top I^{-1}\mu - n\right]$$

$$\therefore KL\left(q(z|x;\theta)||p(z)\right) 與 \mu_\theta(x) \cdot \Sigma_\theta(x) 有关$$