General Idea:

$$z_i \xrightarrow{\hspace{1cm}} \boxed{\theta} \xrightarrow{\hspace{1cm}} x_i$$

$P(z)$     $\underset{\theta}{P(x|z)}$

ML:

$$\max_{\theta} \; \underset{\theta}{P(X)} = \underset{\theta}{P(x_1)} \;-\; \underset{\theta}{P(x_N)}$$

$$\max_{\theta} \; \log P_{\theta}(X) = \sum_{i} \log \underset{\theta}{P(x_i)}$$

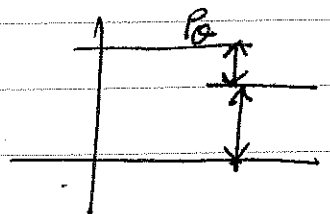$$\log P(x) = \log \int P(x,z)\, dz$$

$$= \log \int P(x|z) P(z)\, dz$$

$$\geq \int \log\left[ P(x|z) P(z) \right] dz$$

$$\frac{1}{2}\left( \log a + \log b \right) \leq \log \frac{a+b}{2}$$

$$\geq \int \log \frac{P_{\theta}(x|z) P(z)}{q(z)} \, q(z)\, dz$$

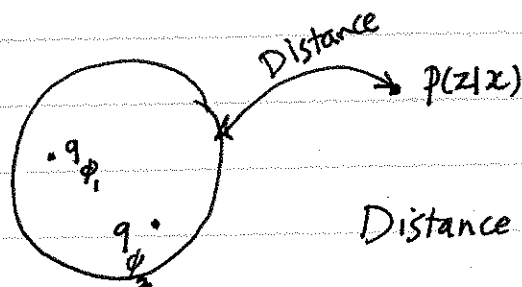$$\underbrace{\qquad\qquad\qquad\qquad}_{\underset{q}{E} \log \frac{P_{\theta}(x,z)}{q_{\psi}(z)}}$$

$$\max_{\theta, \psi} \; E \log \frac{P_{\theta}(x,z)}{q_{\psi}(z)} \qquad (ELBO)$$

Another way of saying the same thing:

$$X = \{x_1, \ldots, x_N\}$$

$$P(Z|X) \quad \longleftarrow \text{ interested}$$

$$P(Z|X) = \frac{P(X,Z)}{P(x)} = \frac{P(x|Z)P(Z)}{P(x)}$$

likelihood

prior

↑ Posterior

$\longleftarrow$ normalizer

$$= \frac{P(x|Z)\,P(Z)}{\int P(x,z)\,dz} \longleftarrow \text{difficult}$$

idea:

$$q_\phi(z)$$

$$\min_\phi \text{Distance}\left( q_\phi(z), p(z|x) \right)$$



Distance $\equiv$ KL

$$KL\left( q_\phi(z), P(z|x) \right) = \mathbb{E}_{q_\phi}\left[ \log \frac{q(z|x)}{q(z)} \right] = \mathbb{E}_{q_\phi}\left[ \log \frac{q_\phi(z)}{P(z|x)} \right]$$

$$= \mathbb{E}_{q_\phi} \log P(z|x)$$

$$= \mathbb{E}_{q_\phi}\left[ \log \frac{q_\phi(z)\,P(x)}{P(x,z)} \right] = \mathbb{E}_{q_\phi} \log \frac{q_\phi(z)}{P(x,z)} + \log P(x)$$

$$\log P(X) = KL\left(q_\phi(z), P(z|X)\right) - \underset{q_\phi}{\mathbb{E}} \log \frac{q_\phi(z)}{P(z,X)}$$

$$\underbrace{\qquad\qquad}_{\ge 0}$$

$$\log P(X) \ge \underset{q_\phi}{\mathbb{E}} \log \frac{P(Z,X)}{q_\phi(z)} \qquad\qquad (ELBO)$$

Expand ELBO:

$$\underset{q}{\mathbb{E}} \log \frac{P(Z,X)}{q_\phi(z)} = \underset{q}{\mathbb{E}} \log \frac{P(X|Z)}{q_\phi(z)} + \underset{q}{\mathbb{E}} \frac{P(Z)}{q_\phi(z)}$$

$$= \underbrace{\underset{q}{\mathbb{E}} \log P(X|Z)}_{\substack{\text{decode} - \text{reconstruction} \\ \text{as if } q \\ \text{is generating}}} - \underbrace{KL\left(q_\phi(z) \| P(z)\right)}_{\text{distance to prior}}$$



$$P(z) \longleftrightarrow q(z) \rightarrow \boxed{\phantom{xx}} \rightarrow X$$
$$\underset{(min)}{\text{distance}} \qquad P(X|Z)$$



$$\underset{q}{\mathbb{E}}\left(\log P(x|z)\right)$$

image $X$ — encode $q(z|x)$ $\phi$ — $Z$ — decode $P(x|z)$

$P(z)$

Challenges:

We have to find best $\phi$ and $\theta$ that max

$$\max_{\phi, \theta} \; -D_{KL}\left(q_\phi(z|x) \| p_\theta(z)\right) + \mathbb{E}_{q_\phi(z|x^i)}\left[\log p_\theta(x^i|z)\right]$$

$$\max_{\phi, \theta} \; -\mathbb{E}_q \log \frac{q(z|x)}{p_\theta(z)} + \mathbb{E}_q\left[\log p_\theta(x^{(i)}|z)\right]$$

$\nabla_\theta$ : easy!

$\nabla_\phi$ : No so much!

Let look closer:

$$\mathbb{E}_{q_\phi}\left[f(z)\right] = \int f(z) \, q_\phi(z) \, dz$$

$$\neq \nabla_\phi \left(\diagup\!\!\!\diagup\right)$$

$$= \int \nabla f(z) q_\phi(z) \, dz = \int f(z) \, \underline{\nabla_\phi q_\phi(z) \, dz}$$

99% you cannot do closed form

$$= \int f(z) \, \frac{\nabla_\phi q_\phi}{q_\phi} \, q_\phi \, dz \qquad \#$$

$$\frac{\nabla_\phi q_\phi}{q_\phi} = \nabla_\phi \log q_\phi$$

Ⓐ $$\int f(z) \nabla \log q_\phi \; q_\phi \, dz = \mathbb{E}\left[ f(z) \nabla \log q_\phi \right]$$

$$\approx \frac{1}{L} \sum f(z_i) \nabla \log q_\phi(z)$$

(we know how to
estimate but it has
large variance!)

Remember ELBO:

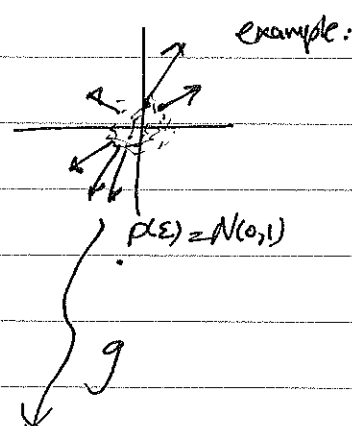$$\mathcal{L} = \int \log \frac{P(x,z)}{q(z|x)} \, q(z|x) \, dz = \mathbb{E}\left[\log P(x,z) - q(z|x)\right]$$

$$\simeq \frac{1}{L} \sum_{\ell} \log\{p(x,z^\ell) - q(z^\ell|x) \qquad \textcircled{I}$$

s.t. $x^\ell \sim q(z|x)$

$\downarrow$ How to generate this?

Sample $\varepsilon^\ell \sim p(\varepsilon)$

$$z^\ell = g(\varepsilon^\ell, x)$$

$\uparrow$
a deterministic function

example:



$p(\varepsilon) = N(0,1)$

$g$

A. Let's look at example of $g$

$$g(\varepsilon) = \begin{bmatrix} 20 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \end{bmatrix}$$
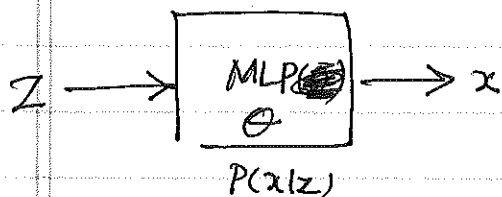$\underset{20}{\overset{?}{\uparrow}}$



Estimation of $\mathbb{E}q\textcircled{I}$ is fine but can have high variance

Let's compute anything that can be computed closed-form!

B. $z = g(\varepsilon) = \mu + \sigma \varepsilon$ - what is $p(z) = ?$
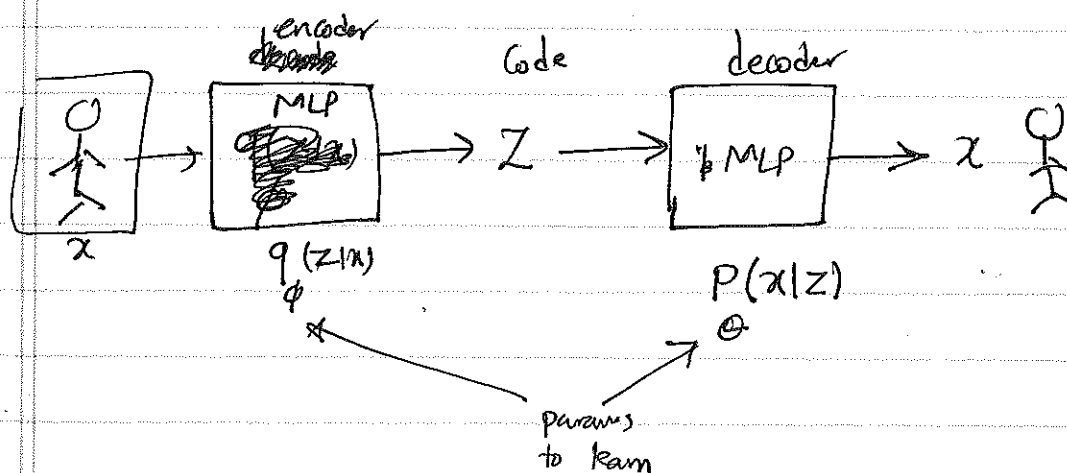
C. $z = g(\varepsilon) = \mu + A\varepsilon$ $p(z) = ?$

our machine



$P(x|z)$

$P(z|x)$   almost impossible
to compute closed-form

Remember   decoder–encoder



encoder

Code    decoder

MLP    $Z$    $\frac{1}{2}$ MLP    $x$

$q_\phi(z|x)$    $P(x|z)$
                  $\theta$

params
to learn

How ~~you~~ we are going to learn $\phi, \theta$ ?

— We set $\phi$ to something, we pass the guy
   to $q_\phi(z|x)$  ¿ how? remember $Z = g_\phi(x, \varepsilon)$ !

   Sample $\varepsilon$, subst. $\varepsilon, x$ get $Z$ !

— Pass $Z$ to decoder get new $\tilde{x}$.

— Is $x, \tilde{x}$ similar? $\mathbb{E}_Z \log P(x|z)$ ↗

~~log~~ ~~log P~~ example $B_?$ $\log P(x|z) = $ ~~$M$~~ $M(x;\theta) + A(x,\theta)z$