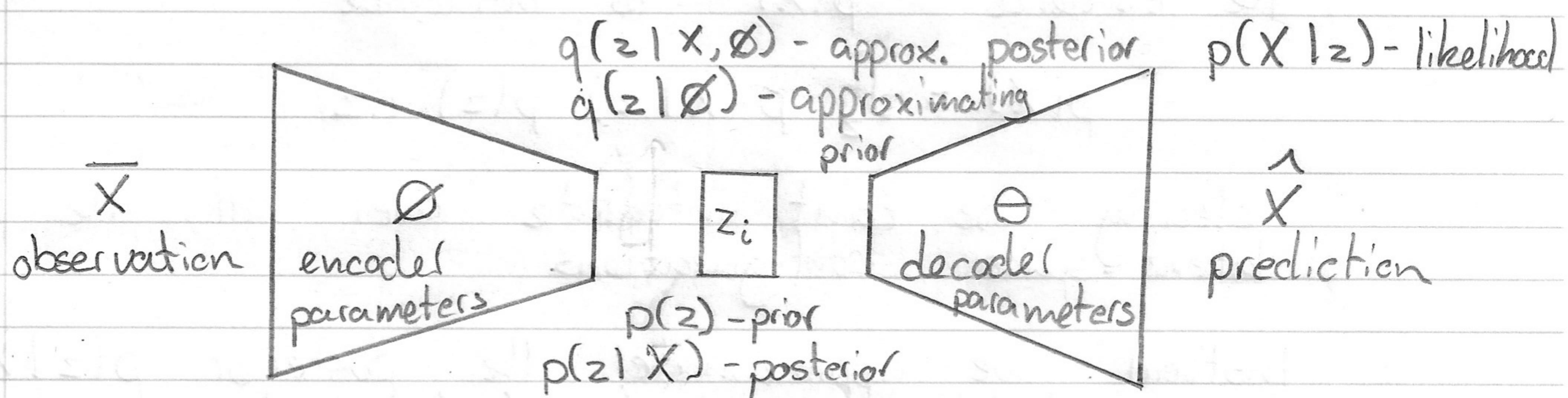


# VAE 1

## Training a VAE



$X = \{\bar{X}, \hat{X}\}$  the set of paired observations or predictions (Note the predictions aren't labels, they are generated!)

we "sample" a distribution  $z_i \sim p(z)$   
please look this up if you don't understand it!

It is not possible to take derivatives using SGD to train this, we will see why not mathematically

We can determine the posterior  $p(z|\bar{X})$  using Bayes' equation

$$p(z|\bar{X}) = \frac{p(\bar{X}|z)p(z)}{p(\bar{X})}$$

where  $p(\cdot)$  is a probability distribution

for example the likelihood can be a gaussian

$$p(\bar{X}|z) = N(0, \Sigma)$$

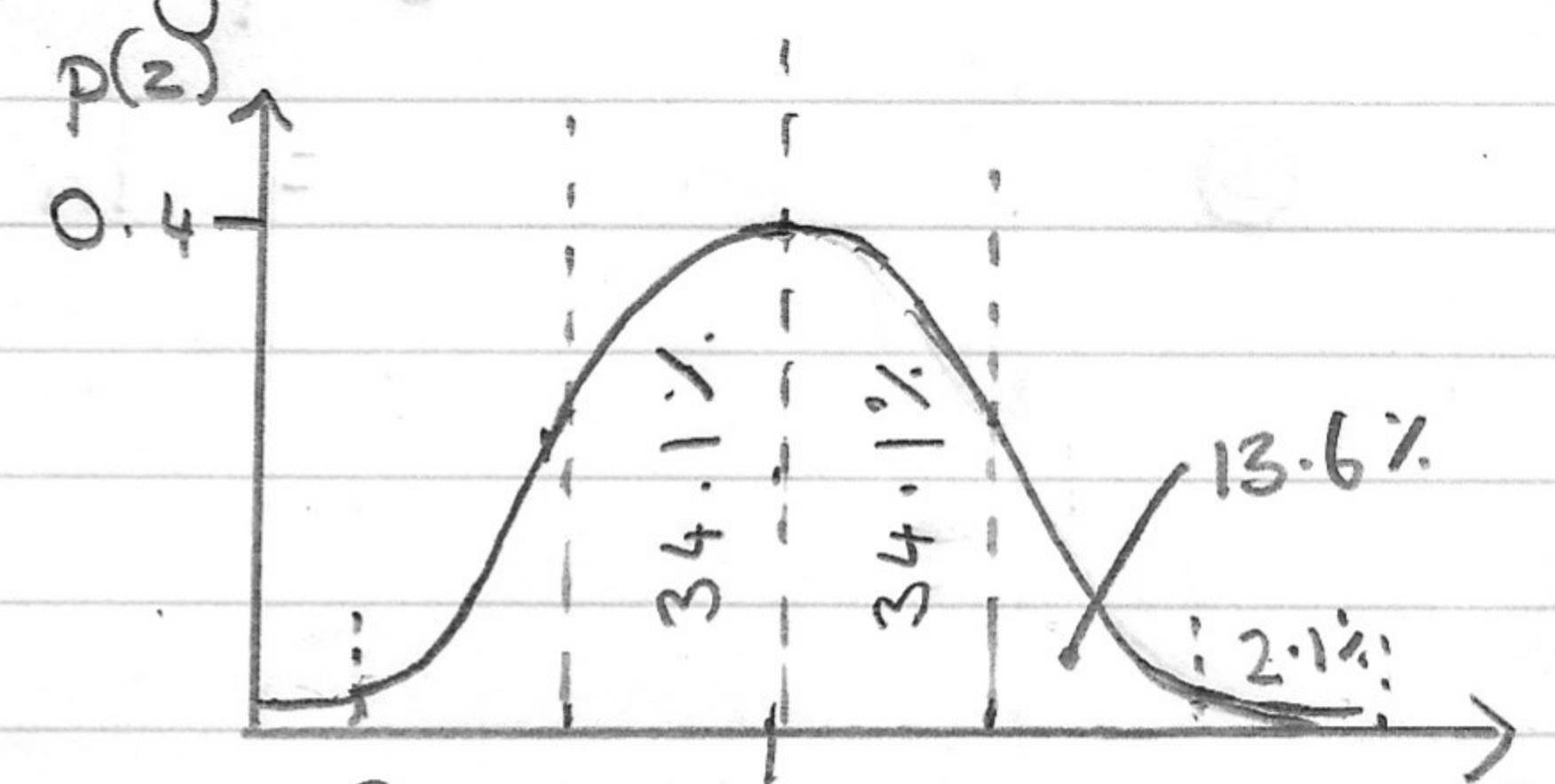
$$= \frac{1}{\sqrt{2\pi} \Sigma} \exp\left(-\frac{1}{2} \frac{\bar{X}^2}{\Sigma^2}\right)$$

where  $\bar{X} = \bar{X} - \hat{X}$

or the prior can be a gaussian

$$p(z) = N(0, 1)$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z^2\right)$$



where 95.4% values between  $[-2, 2]$   
99.7% values between  $[-3, 3]$

We can't solve Bayes' equation directly because the evidence  $p(X)$  is intractable

$$p(X) = \int_z p(X|z) p(z) dz$$

clearly we can't integrate over all the possible latent-space configurations.

Instead we approximate the posterior  $p(z|X)$  with another distribution  $q(z|X, \emptyset)$  which depends on the output of the encoder.

We find the best approximating distribution when the KL divergence is minimised

$$\textcircled{A} \quad q(X|z, \emptyset^*) = \min_{\emptyset} \text{KL}[q(z|X, \emptyset) \| p(z|X)]$$

Clearly this expression still has issues because it contains the posterior

### ELBO

The Evidence Lower Bound is a way to rearrange Eq A so we can solve it

$$\begin{aligned} \text{KL}[q_{\emptyset} \| p(z|X)] &= \int_z q_{\emptyset} \ln \left[ \frac{q_{\emptyset}}{p(z|X)} \right] dz \\ &= \int_z q_{\emptyset} \ln \left[ \frac{q_{\emptyset} p(X)}{p(z, X)} \right] dz \\ &= \int_z q_{\emptyset} \ln[p(X)] + q_{\emptyset} \ln \left[ \frac{q_{\emptyset}}{p(z, X)} \right] dz \end{aligned}$$

Note  $\int_z q(z|X, \emptyset) dz = 1$

$$\textcircled{B} \quad = \ln[p(X)] + \int_z q_{\emptyset} \ln \left[ \frac{q_{\emptyset}}{p(z, X)} \right] dz$$

We also know that  $KL[q_\phi \parallel p(z|x)] \geq 0$  by definition, so we can rewrite Eq B

$$0 \leq \ln[p(x)] + \int_z q_\phi \ln \left[ \frac{q_\phi}{p(z,x)} \right] dz$$

$$\textcircled{C} \quad \ln[p(x)] \geq -KL[q_\phi \parallel p(z,x)]$$

So the evidence is bounded below by the negative of the KL divergence of the approximating & joint distributions

REMEMBER probability distributions are always less than 1 so

$$p(x) \leq 1 \quad \text{therefore} \quad \ln[p(x)] \leq 0$$

(Eq A) It should be possible to see that minimising  $KL[q_\phi \parallel p(z|x)]$  involves minimising the <sup>log</sup> evidence, which can be found by minimising the Right Hand Side of Eq C.

We want to use a gradient based optimiser, but how do we take a gradient of Eq C?

### Gradient of the KL divergence

We can rewrite Eq A using Eq C.

$$\begin{aligned} q(z|x, \phi^*) &= \min_{\phi} -KL[q_\phi \parallel p(z,x)] \\ \textcircled{D} \quad &= \min_{\phi, \theta} \underbrace{\int_z q_\phi \ln [p(x|z, \theta, \phi)]}_{\text{Term 1}} - \underbrace{q_\phi \ln \left[ \frac{q_\phi}{p(z)} \right]}_{\text{Term 2}} dz \end{aligned}$$

When we perform SGD we want the gradient estimate at each iteration to be correct "on average"

The mathematical way to do this is say that "the gradient of the expectation of  $p(x|z, \theta, \phi)$ " should be the same as "the expectation of the gradient of  $p(x|z, \theta, \phi)$ ".

Unfortunately we can't do this transform directly as it leads to a "Biased estimator".

Note the definition of the expectation of  $f(y, x)$  with respect to a probability  $p(x)$

$$E_{p(x)}[f(y, x)] = \int_x p(x) f(y, x) dx$$

So for Term 1 in Eq D

$$\begin{aligned} \nabla_\theta E_{q_\theta}[p(x|z, \theta, \phi)] &= \nabla_\theta \int_z q(z|x, \phi) \ln[p(x|z, \theta, \phi)] dz \\ &\neq \int_z q(z|x, \phi) \nabla_\theta \ln[p(x|z, \theta, \phi)] dz \\ &\neq E_{q_\theta}[\nabla_\theta \ln[p(x|z, \theta, \phi)]] \end{aligned}$$

This means when we have a specific observation  $X_i \sim p(X)$  we can't use Monte Carlo methods to take the derivative wrt.  $\theta$

wrt  
"with  
respect to"  
How do we get to an "unbiased estimator of the gradient"?

### The log derivative trick

The "log derivative trick" or "score function method" is the first of two approaches to finding the unbiased gradient estimator

$$\begin{aligned} \nabla_\theta E_{q_\theta}[\ln[p(x|z, \theta, \phi)]] &= \int_z \nabla_\theta [\ln[p(x|z, \theta, \phi)]] q(z|x, \phi) dz \\ &= \int_z q(z|x, \phi) \nabla_\theta \ln[p(x|z, \theta, \phi)] + \ln[p(x|z, \theta, \phi)] \nabla_\theta q(z|x, \phi) dz \\ &= \int_z q_\theta \nabla_\theta \ln[p(x|z, \theta, \phi)] + \ln[p(x|z, \theta, \phi)] q_\theta \nabla_\theta \ln[q_\theta] dz \\ &= E_{q_\theta}[\nabla_\theta \ln[p(x|z, \theta, \phi)] + p(x|z, \theta, \phi) \nabla_\theta \ln[q_\theta]] \end{aligned}$$

where we use the rule  $\nabla_x f(x) = f(x) \nabla_x \ln[f(x)]$

Great!

The problem is that this estimator has a high variance because  $0 \leq q_\theta \leq 1$  which means  $\ln[q_\theta] \ll 0$  (sometimes)

This is a problem once we sample to take the gradients

We have as Eq F

$$z_i \sim q(z | X, \emptyset)$$

$$\hat{x}_i \sim p(x | z_i, \theta, \emptyset)$$

$$\nabla_{\theta} \mathbb{E}_{q(\cdot)}[p(x | z, \theta, \emptyset)] = n \left[ \nabla_{\theta} (\bar{x}_i - \hat{x}_i)^2 + (\bar{x}_i - \hat{x}_i)^2 \nabla_{\theta} z_i^2 \right]$$

which is the MC gradient estimate,  $\bar{x}_i$  for example is a single observed image

### The reparameterisation trick

This is an approach which reduces the variance of the gradient estimates

Sometimes the distribution chosen for the latent variable can be rewritten as a "deterministic function" of another distribution with constant parameters

Mathematically  $q(z | X, \emptyset) = g(p(\epsilon), \{\lambda\})$ .

An example is necessary to understand this. Consider a normally distributed latent space

$$q(z | X, \emptyset) = N(\mu, \Sigma)$$

In general normal distributions

$$g(p(\epsilon), \{\lambda\}) = \mu + \Sigma N(0, 1)$$

$$\text{where } p(\epsilon) = N(0, 1)$$

$$\{\epsilon\} = \{0, 1\} \text{ NOTE these are constants!}$$

$$\{\lambda\} = \{\mu, \Sigma\} \text{ (the parameters of the gaussian)}$$

This means we now have differentiable parameters in the latent space, where the parameters are each functions of the encoder  $\{\lambda\} = \{\mu(\emptyset), \Sigma(\emptyset)\}$

This means Term 1 Eq D becomes

$$q(z | X, \emptyset^*) = \min_{\emptyset, \theta} \int_z q(z | X, \emptyset) \ln [p(x | z, \theta, \emptyset)] dz$$

$$= \min_{\emptyset, \theta} \int_z p(\epsilon | g(N(0, 1), \mu, \Sigma)) \ln [p(x | z, \theta, \emptyset)] dz$$

Differentiating this expression is really straight-forward

Note

$$\begin{aligned} q(\epsilon) &= N(0, 1) \quad \nabla_{\theta} E_{q(\epsilon)} [g(N(0, 1), \mu, \Sigma) \ln[p(x|z, \theta, \emptyset)]] \\ &= \int_{\epsilon} q(\epsilon) \nabla_{\theta} [g(q(\epsilon), \mu, \Sigma) \ln[p(x|z, \theta, \emptyset)]] d\epsilon \\ &= \int_{\epsilon} q(\epsilon) \left[ \ln[p(x|z, \theta, \emptyset)] \underline{\nabla}_{\mu} g(q(\epsilon), \mu, \Sigma) \nabla_{\theta} \mu(\emptyset) + \dots \right. \\ &\quad \left. \dots g(q(\epsilon), \mu, \Sigma) \nabla_{\theta} \ln[p(x|z, \theta, \emptyset)] \right] d\epsilon \end{aligned}$$

but does look a bit complicated

Note that an additional expression is required for each additional parameter of  $g(\cdot)$ . In this case we will have one additional expression where  $\underline{\nabla}_{\mu}$  will be replaced by  $\underline{\nabla}_{\Sigma}$