

AE + VAE - Rough Overview

- The principal aim of an autoencoder is to find a meaningful low-dimensional (relatively speaking) representation of data
- e.g. PCA could be viewed as a naive autoencoder, capturing linear features.
- Given a density γ , denote by γ^N the empirical density,
$$\gamma^N(u) := \frac{1}{N} \sum_{n=1}^N f(u - u^{(n)})$$

↑
sampled point.
- The aim is to find an approximate factorisation of Id which is accurate on the support of γ , given access only to γ^N .

net $f : \mathbb{R}^d \times \Theta_f \rightarrow \mathbb{R}^{d_2}$
 $g : \mathbb{R}^{d_2} \times \Theta_g \rightarrow \mathbb{R}^d$.

e.g., f and g can be NNs with parameters w_{ij} in the spaces $\Theta_f \subseteq \mathbb{R}^{P_f}$ and $\Theta_g \subseteq \mathbb{R}^{P_g}$, respectively

- Then simply define the reconstruction loss

$$J(\Theta_f, \Theta_g) = \| \text{Id} - g(f(\cdot; \Theta_f) \Theta_g) \|_{\text{H}_Y}^2$$

$$:= \mathbb{E}_{u \sim \gamma} [\| u - g(f(u; \Theta_f) \Theta_g) \|^2]$$

Obviously we have to compute this expectation w.r.t. γ^N in reality.

We thus have the objective

$$(\Theta_f^*, \Theta_g^*) \in \arg \min_{(\Theta_f, \Theta_g)} J(\Theta_f, \Theta_g),$$

and write $f^*(u) = f(u; \theta_f^*)$ and $g^*(u) = g(u; \theta_g^*)$

- Note that we have converted an unsupervised learning problem into a supervised learning problem.
- Some benefits include robustness, computational tractability, and interpretability.

PCA Example

We elaborate a bit on the 'PCA as an autoencoder' comment.

Let $m := \mathbb{E}[u]$ under $u \sim \gamma^n$

and $C := \mathbb{E}[(u-m) \otimes (u-m)]$ under $u \sim \gamma^n$.

(\otimes denotes the kronecker/tensor product, so
 $(C)_{ij} = \mathbb{E}[(u_i - m_i)(u_j - m_j)] = (C)_{ji}$)

PCA consists of finding at its first step

$$w_1 = \arg \max_{\|w\|=1} \left\{ \sum_i \langle u_i, w \rangle^2 \right\}$$

$$= \arg \max \left\{ \frac{w^T U^T U w}{w^T w} \right\},$$

where U is the mean-shifted data matrix.

i.e. $U^T U = C$, the covariance matrix.

This is maximised when w is the eigenvector corresponding to the largest eigenvalue. The point is to transfer to a coordinate system where the greatest variance of the data among all subspace projections is given by that spanned by $\{w_i\}_{i=1}^k$, for k the desired # of components.

To get further components one subtracts the first component from the data matrix, $U_2 = U - U w_1 w_1^T$, and repeat the process.

So, let $(w_i^i, \gamma_i^i) \in \mathbb{R}^d \times \mathbb{R}^+$ be the eigenpairs (normalised) s.t. $\gamma_1 > \gamma_2 > \dots > \gamma_n$. Note that C is symmetric and PSD, so all eigenvalues are positive and there exists an orthonormal basis of eigenvectors.

For $d_2 \leq d$, define

$$f(u) = (\langle u, w_1 \rangle, \dots, \langle u, w_{d_2} \rangle) \in \mathbb{R}^{d_2}$$

$$g(z) = \sum_{j=1}^{d_2} z_j w_j$$

$$\rightarrow g(f(u)) = \left(\sum_{j=1}^{d_2} w_j \otimes w_j \right) u$$

i.e., this is a truncated version of the representation of Id given by the spectral theorem for PD symmetric matrices.

VAEs

- Define the pushforward map by

$$f_{\#} \mu(B) := \mu(f^{-1}(B)), \quad B \in \mathcal{E}_2,$$

where $(X_1, \Sigma_1), (X_2, \Sigma_2)$ are measurable spaces, and f is μ -measurable map from X_1 to X_2 , and

$$\mu: \Sigma_1 \rightarrow [0, +\infty)$$

- Then

$$\mathbb{f} = f^* \# \gamma$$

gives the approximate distribution in the latent space, but

we can't specify ξ .

- we might want to impose some specified distribution in the latent space. This is the point of a VAE.

Definition:

- * Given densities ρ and ℓ on \mathbb{R}^d ,
a coupling is a density π on $\mathbb{R}^d \times \mathbb{R}^d$
s.t.

$$\int_{\mathbb{R}^d} \pi(z, u) du = \rho(z)$$

and

$$\int_{\mathbb{R}^d} \pi(z, u) dz = \ell(u)$$

so marginals of deliver ρ and ℓ .

- Note that couplings are by no means unique.
- The extension to ρ and ℓ means one different dimension is obvious.

- Let π be a coupling of γ and ξ . Then
by definition, we have

$$\pi(u, z) = \underbrace{\text{IP}(u|z)}_{\curvearrowleft} \xi(z) = \underbrace{\text{IP}(z|u)}_{\curvearrowright} \gamma(u)$$

The point of a VAE is to
approximate the conditional densities
by parametrised families.

- for example, we could have

$$u|z \sim N(g(z; \theta_g), \sigma_g^2 \text{Id})$$

$$z|u \sim N(f(u; \theta_f), \sigma_f^2 \text{Id})$$

The approximate coupling density becomes

the different approximations of the same thing!

$$\begin{aligned} n_g(u, z) &= \frac{1}{Z_g} \exp\left(-\frac{1}{2\sigma_g^2} |u - g(z; \theta_g)|^2\right) \zeta(z) \\ n_f(u, z) &= \frac{1}{Z_f} \exp\left(-\frac{1}{2\sigma_f^2} |z - f(u; \theta_f)|^2\right) \gamma(u) \end{aligned}$$

for simplicity, let's assume σ_f and σ_g are fixed and determine the parameters of f and g , meaning find θ_f and θ_g s.t. n_g and n_f are closed.

Definition:

* The KL-divergence of a distribution P from a distribution Q is

$$D_{KL}(P || Q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

Lemma:

* If $p = N(\mu_p, \Sigma_p)$ and $q = N(\mu_q, \Sigma_q)$, then

$$\begin{aligned} D_{KL}(p || q) &= \frac{1}{2} \left[\text{Tr}(\Sigma_q^{-1} \Sigma_p) - d + (\mu_p - \mu_q)^T \Sigma_q^{-1} (\mu_p - \mu_q) \right. \\ &\quad \left. + \log\left(\frac{\det \Sigma_q}{\det \Sigma_p}\right) \right] \end{aligned}$$

One can check that $\downarrow \sim N(f(u; \theta_f), \sigma_f^2 \text{Id})$

$$D_{KL}(n_f || n_g) = E_{u \sim \gamma} [E^{z \sim P(z|u)} \left[\frac{1}{2\sigma_g^2} |u - g(z; \theta_g)|^2 \right]]$$

$$+ D_{KL}(P(z|u) || \zeta) + \text{const}$$

$$= \frac{1}{2} \|f(u; \theta_f)\|^2 \quad \text{by}$$

$$= E^{u \sim r} [E^{\xi \sim N(0,1)} [\frac{1}{2\sigma_g^2} |u - g(f(u; \theta_f) + \tau_f \xi; \theta_g)|^2] \\ + \frac{1}{2} \|f(u; \theta_f)\|^2]$$

This is a regularised version of the standard autoencoder.

VAEs in greater detail

- Assume that $x \sim p^*$, where we will print write the density as $p^*(x)$. This is a true distribution we need to model.
- we want a model class which is flexible enough to adapt to data whilst respecting a priori constraints. we denote our model by $p_\theta(x)$
- Denote DNNs as a neural function Neural Net (\circ). we will be interested in using DNNs to approximate density functions.

Latent Variables

- latent variables, z , are variables which we do not see, so we cannot infer their values directly.
 - Our data model then represents a joint distribution, $p_\theta(x, z)$, over both observed and latent variables
 - The marginal distribution is just
- $$p_\theta(x) = \int p_\theta(x, z) dz,$$
- also called the (single datapoint) marginal likelihood (model evidence).
- Deep latent Variable Models are latent variable models for which $p_\theta(x, z)$ has distributions parametrized by NNs.
 - The benefit is that we can have a simple directed model, like a parametrised Gaussian conditional, but build them up to be very expressive.
 - Formally, this means that

$$p_\theta(x, z) = p_\theta(z) p_\theta(x|z),$$

where one or both of the densities on the RHP is specified. $p(z)$ is the prior.

- The difficulty in training these models is that for DLMs the marginal probability of data under the model is hard to estimate.
- This means we can't differentiate w.r.t. θ and thus optimise.
- VAEs are designed to address this problem.

Variational Autoencoders

- Note that $p_\theta(x, z)$ is tractable, and that

$$p_\theta(z|x) = \frac{p_\theta(x, z)}{p_\theta(x)}.$$

So to get $p_\theta(x)$, can we find a good way of determining roughly what $p_\theta(z|x)$ is?

- Yes, introduce a parametric inference model, $q_\phi(z|x)$.
This is called the encoder.
- Similarly to a DLM, we can parametrise $q_\phi(z|x)$ by a NN. The variational parameters ϕ are the weights and biases of the NN.

e.g. $(\mu, \log(\sigma)) = \text{Encoder Neural Net}_\phi(x)$

$$q_\phi(z|x) = \mathcal{N}(z; \mu, \sigma)$$

- So we see that for each point x we're trying to learn the parameters which give a informative distribution in the latent space.

The Evidence Lower Bound

- Most variational Bayesian methods try to maximize a quantity called the ELBO
- for any choice of inference model, $q_\phi(z|x)$, we have

$$\log(p_\theta(x)) = E_{z \sim q_\phi(z|x)} [\log(p_\theta(x))]$$

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

$$= E_{z \sim q_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{q_\phi(z|x)} \cdot \frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right]$$

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

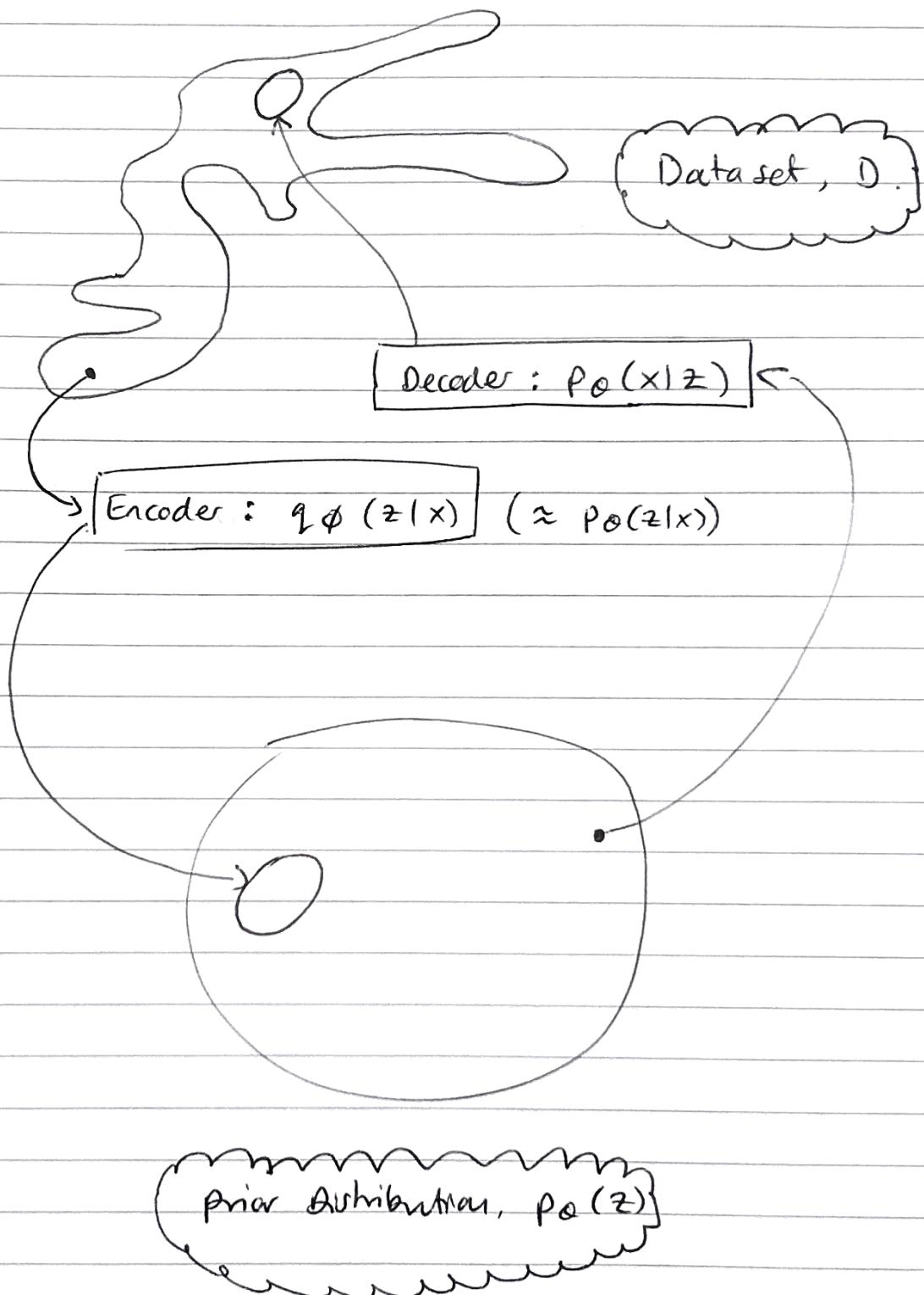
$$+ E_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right]$$

$$:= \underbrace{L_{\theta, \phi}(x)}_{\text{ELBO}(x)} + \underbrace{D_{KL}(q_\phi(z|x) || p_\theta(z|x))}_{\geq 0},$$

so this is indeed a lower bound on the probability of that data x under the model p_θ .

- Note that we have equality iff $q_\phi(z|x)$ equals the true posterior distribution.
- So in addition to, by definition, quantifying the divergence of the approximate posterior from the true posterior, we can interpret D_{KL} as the gap between the ELBO and the marginal likelihood $p_\theta(x)$.

To summarise schematically, we are doing the following:



So we see that the generative model learns a joint distribution $p_\theta(x, z)$ factorised as $p_\theta(z)p_\theta(x|z)$, where we impose a prior $p_\theta(z)$ on the latent space.

The stochastic encoder $q_\phi(z|x)$ approximates the true but intractable posterior $p_\theta(z|x)$ of the model.

- we just use the ELBO, $\mathbb{E}[\log(p_\theta(x, z)) - \log(q_\phi(z|x))]$ under $z \sim q_\phi(z|x)$, as the optimisation objective, jointly finding the best parameters Θ and ϕ . e.g. using SGD.
- As our previous note on the interpretation of the ELBO said, this optimisation gives us :
 - Approximate maximisation of the marginal likelihood, $p_\theta(x)$, so our generative model improves
 - Minimisation of the KL divergence of $q_\phi(z|x)$ from the true posterior, $p_\theta(z|x)$, making $q_\phi(z|x)$ better.

Estimating the gradient

- The individual datapoint ELBO, as well as its gradient, $\nabla_{\Theta, \phi} \mathcal{L}_{\Theta, \phi}(x)$, are normally intractable, but good estimators exist.
- for gradients w.r.t. the parameters of the generative model, write

$$\begin{aligned}\nabla_{\Theta} \mathcal{L}_{\Theta, \phi}(x) &= \nabla_{\Theta} \mathbb{E}_{z \sim q_\phi(z|x)} [\log(p_\theta(x, z)) - \log(q_\phi(z|x))] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\nabla_{\Theta} (\log(p_\theta(x, z)) - \log(q_\phi(z|x)))] \\ &\simeq \nabla_{\Theta} \log(p_\theta(x, z)),\end{aligned}$$

where we use \simeq to denote that this is at the end w.r.t. a random sample of z . i.e., use Monte Carlo to estimate gradients.

- Gradients w.r.t. ϕ are more challenging, since the expectation is taken w.r.t. $q_\phi(z|x)$ which depends on ϕ . i.e., $\nabla_{\phi} \mathbb{E}_{z \sim q_\phi(z|x)} [\log(p_\theta(x, z)) - \log(q_\phi(z|x))]$ is not generally

the same as $E_{z \sim q_\phi(z|x)}[-\nabla_\phi \log(q_\phi(z|x))]$.

- To deal with this problem, we introduce a new parametrization

parametrization Trick

- write $z \sim q_\phi(z|x)$ as the transformation of some other RV ϵ , given x and ϕ :

$$z = g(\epsilon, \phi, x),$$

where the distribution of ϵ is independent of x and ϕ .

- note that g must be invertible and differentiable.

- Now we can write

$$\begin{aligned}\nabla_\phi E_{z \sim q_\phi(z|x)}[f(z)] &= \nabla_\phi E_{\epsilon \sim p(\epsilon)}[f(z)] \\ &= E_{\epsilon \sim p(\epsilon)}[\nabla_\phi f(z)] \\ &\simeq \nabla_\phi f(z),\end{aligned}$$

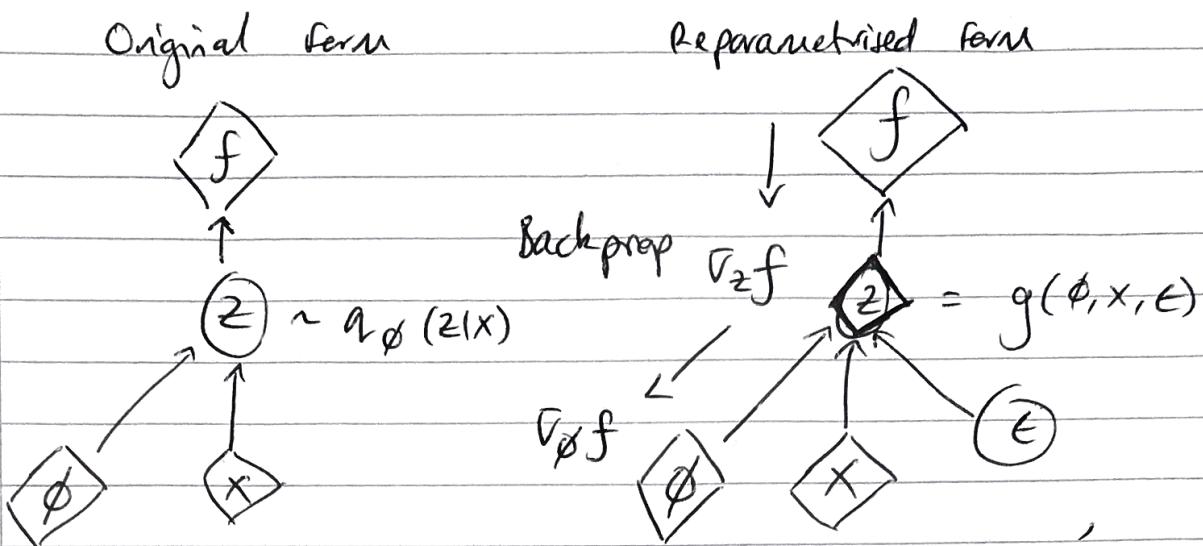
for $z = g(\phi, x, \epsilon)$ and $\epsilon \sim p(\epsilon)$ a random sample.

- One can check that this single sample algorithm, in which we:

- a) draw $\epsilon \sim p(\epsilon)$
 - b) set $z = g(\phi, x, \epsilon)$
 - c) estimate $\tilde{\mathcal{L}}_{\phi, \epsilon}(x) = \log(p_\phi(x, z)) - \log(q_\phi(z|x))$
- is an unbiased estimator of the gradient:

$$\begin{aligned}E_{p(\epsilon)}[\nabla_{\phi, \epsilon} \tilde{\mathcal{L}}_{\phi, \epsilon}(x; \epsilon)] &= E_{z \sim p(\epsilon)}[\nabla_{\phi, \epsilon} \log(p_\phi(x, z)) - \nabla_{\phi, \epsilon} \log(q_\phi(z|x))] \\ &= \nabla_{\phi, \epsilon} E_{z \sim p(\epsilon)}[\log(p_\phi(x, z)) - \log(q_\phi(z|x))] \\ &= \nabla_{\phi, \epsilon} \mathcal{L}_{\phi, \epsilon}(x).\end{aligned}$$

Schematically, we have the following conversion:



where $\diamond = \text{deterministic}$ and $\circ = \text{random}$.

Computation of the ELBO

- To evaluate the estimator of the ELBO, we need the density, $\log(q_\phi(z|x))$. This is easy if we choose the right $g(\cdot)$, noting that typically one knows (chooses, indeed) $p(\epsilon)$.
- If g is invertible, the densities of ϵ and z are related via

$$\log(q_\phi(z|x)) = \log(p(\epsilon)) - \log(D_\phi(x, \epsilon)),$$

$$\text{for } D_\phi(x, \epsilon) = \frac{\partial z}{\partial \epsilon} = \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \dots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \dots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix},$$

The Jacobian of the transformation $\epsilon \mapsto z(\epsilon)$.

- The hope is that one can build transformations, $g(\cdot)$, which are flexible and for which $\log(D_\phi(x, \epsilon))$ is easy to compute (hence giving a tractable and flexible

inference model, $q_{\phi}(z|x)$)

An Example

- we will treat the case of a Gaussian with diagonal covariance matrix, and then the full case.

① we will say that

$$q_{\phi}(z|x) = N(z; \mu, \text{diag}(\sigma^2)),$$

where μ and σ^2 are functions of x :

$$(\mu, \log(\sigma)) = \text{Encoder Neural Net } (x)$$

$$\rightarrow q_{\phi}(z|x) = \prod_i q_{\phi}(z_i|x) = \prod_i N(z_i; \mu_i, \sigma_i^2)$$

for the reparametrisation, write

$$\epsilon \sim N(0, \text{Id})$$

$$z = \mu + \sigma \odot \epsilon, \quad (= g(\epsilon, \phi, x))$$

for \odot the element-wise product.

obviously the Jacobian is $\text{diag}(\sigma)$, so the 'log determinant' is $\sum_i \log(\sigma_i)$

$$\Rightarrow \log(q_{\phi}(z|x)) = \log(p(\epsilon)) - \log(\det(\text{J}\phi(x, \epsilon)))$$

$$= \sum_i \log(N(\epsilon_i; 0, 1)) - \log(\sigma_i)$$

② for the full-covariance Gaussian posterior,

$$q_{\phi}(z|x) = N(z; \mu, \Sigma)$$

The appropriate parametrisation is

$$\epsilon \sim N(0, \text{Id})$$

$$z = \mu + L\epsilon$$

for L some lower triangular matrix. We can parametrise this way since given A s.t. $AA^T = \Sigma$, a normal RV is $z = \mu + A\epsilon$.

But now A is not unique, and for any symmetric positive-definite matrix we can use the Cholesky decomposition to find a lower (or upper) triangular A .

i.e. L is

$$\begin{aligned}\Sigma &= E[(z - E[z])(z - E[z])^T] \\ &= E[(L\epsilon)(L\epsilon)^T] \\ &= L E[\epsilon\epsilon^T] L^T \\ &= LL^T.\end{aligned}$$

We can build such an L by having

$$(\mu, \log(\tau), L') = \text{Encoder Neural Net } (x)$$

$$L = L_{\text{mask}} \odot L + \text{diag}(\tau),$$

where $L_{\text{mask}} = \begin{pmatrix} 0 & \cdots & 0 \\ 1 & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & 0 \\ 1 & \cdots & 0 & 0 \end{pmatrix}$, so 0 on and above the main diagonal.

so as before, the log-determinant is $\sum_i \log(\tau_i)$.

- One can see that, in general, if $z = L\epsilon + \mu$ can be replaced with a chain of potentially non-linear transformations and provided the Jacobian is triangular with non-zero diagonal entries the log-determinant remains easy.

Estimation of the Marginal likelihood

- After training a VAE, we can estimate the probability of data under the model. Obviously we have

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

$$\approx \log \left(\frac{1}{L} \sum_{e=1}^L p_\phi(x, z^e) / q_\phi(z^e|x) \right)$$

for z^e i.i.d. $z^e \sim q_\phi(z|x)$

Generative Modelling

- We can now obviously draw random or conditional new samples from an learned distribution.
- For instance, to generate samples similar to x , do the following:

(1) Find the parameters $\phi = \text{Neural Net}(x)$ which parametrise q_ϕ .

i.e. use the encoder to get
 $q_\phi(z|x)$

(2) Sample z from this distribution.

(3) Decode z to get ~~$x_{\text{new}} = \text{Decoder Neural Net}(z)$~~
 $x_{\text{new}} = \text{Decoder Neural Net}(z)$

- Note that there is some subtlety in (3). Strictly, we learn parameters of $p(x|z)$ which outputs a mean function.

- We don't normally have a variance estimate here. This is because selecting the best point minimises reconstruction loss, in that if $p_\phi(x|z)$ is $N(\mu_\phi, \sigma_\phi^2 \text{Id})$,

then $E_{q(z|x)} [-\log(p_0(x|z))]$ is ~~maximized~~ maximized by ~~approximation~~ minimizing the squared error of x and $M_\phi(z)$.

- Note that really we could view ϕ as the encoding parameters and θ as the decoding (generative modelling) parameters, so far unconditional modelling we don't need ϕ at all! (Although obviously we still need it in order to learn p_0).