

Lecture notes on variational auto-encoders (VAEs)

1.1 Prerequisites

Divergence minimisation

Given a set of probability distributions \mathcal{P} on a random variable X , a divergence is defined as a function $D[\cdot||\cdot] : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$ such that $D[P||Q] \geq 0$ for all $P, Q \in \mathcal{P}$, and $D[P||Q] = 0$ iff. $P = Q$.

The definition of divergence is much weaker than that for a *distance* such as the ℓ_2 -norm, since it does not need to satisfy either symmetry in arguments or the triangle inequality. There exist many available divergences to use, some of them will be introduced throughout this course.

In this course we assume the probability distributions/measures in \mathcal{P} are dominated by the Lebesgue measure of the underlying Euclidean space, so that we can work with *probability density functions* (PDFs)

$$p(\mathbf{x}) = \frac{dP}{d\mathbf{x}}, \forall P \in \mathcal{P}. \quad (1)$$

In the following we will also write \mathcal{P} as the set of PDFs w.l.o.g., and use the terms probability distribution and probability density functions interchangeably (unless specifically mentioned).

Probabilistic graphical models

In machine learning tasks we may define the model as a distribution on a set of random variables with particular dependency structures. Some of the variables might be unobserved as well. Probabilistic graphical models are powerful models that use graphs to describe the dependency structure of the random variables. In particular we consider direct acyclic graphs (DAGs) which are graphs with directed edges and without directed cycles. By assuming Markov properties, DAGs can be used to describe the factorisation structure of the joint distribution. Interested readers are referred to e.g. Chapter 8 of Bishop [2007] for a formal introduction of probabilistic graphical models. For this course we only introduce the principles for reading joint distributions from a DAG (and vice versa). Assuming we are interested in the distribution $p(\mathbf{x}_1, \dots, \mathbf{x}_D)$ for a given DAG with nodes $\{\mathbf{x}_1, \dots, \mathbf{x}_D\}$ and directed edges between them, then the joint distribution is:

$$p(\mathbf{x}_1, \dots, \mathbf{x}_D) = \prod_{i=1}^D p(\mathbf{x}_i | pa(\mathbf{x}_i)), \quad (2)$$

where $pa(\mathbf{x}_i) \subset \{\mathbf{x}_1, \dots, \mathbf{x}_D\}$ represents the parent nodes of \mathbf{x}_i in the DAG. For a DAG there always exists root node(s) that have no parents (i.e. $pa(\mathbf{x}_i) = \emptyset$), and in such case $p(\mathbf{x}_i | pa(\mathbf{x}_i)) = p(\mathbf{x}_i)$. Conversely, given a joint distribution in the form of (2), we can also draw the corresponding DAG by adding arrows from nodes in $pa(\mathbf{x}_i)$ to \mathbf{x}_i . A number of examples are visualised in Figure 1.

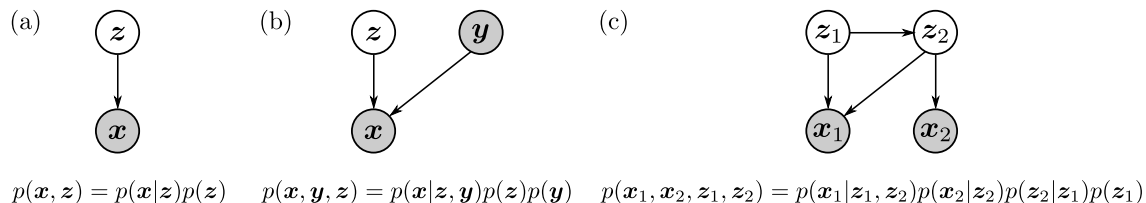


Figure 1: Examples of probabilistic graphical models: graphs & the corresponding factorisations of the joint distributions. Shaded nodes represent observed variables and the other nodes represent unobserved/latent variables. Example (a) corresponds to the latent variable model used in VAEs & GANs, and example (b) corresponds to the latent variable model used in conditional VAEs & GANs where \mathbf{y} represents additional information that the generative model is conditioned on.

Jensen's inequality

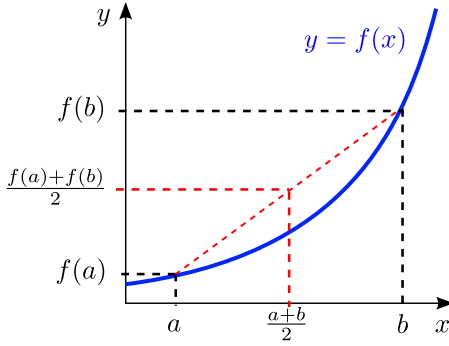
Below we introduce Jensen's inequality as a prerequisite for later discussions on divergences.

Proposition 1. (*Jensen's inequality*) If $f : \mathbb{R} \rightarrow \mathbb{R}$ is a convex function, then for any distribution $p(x)$,

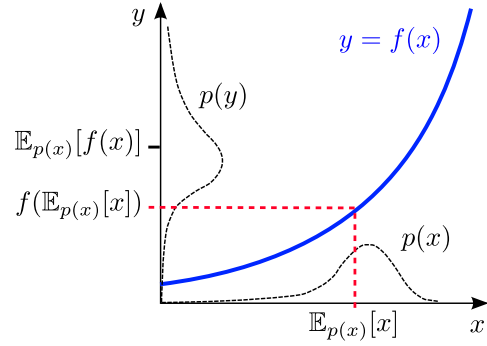
$$\mathbb{E}_{p(x)}[f(x)] \geq f(\mathbb{E}_{p(x)}[x]),$$

with equality holds iff. f is linear or $p(x)$ is a delta measure.

A visual proof is provided in the below figures.



(a) for discrete distributions



(b) for continuous distributions

Jensen's inequality can be generalised to functions formed by compositions of functions. To see this, we first introduce the *law of the unconscious statisticians* (LOTUS) rule:

Proposition 2. (*LOTUS*) Given a distribution $p_X(\mathbf{x})$ and a function $\mathbf{y} = g(\mathbf{x})$ such that $\mathbb{E}_{p_X(\mathbf{x})}[g(\mathbf{x})] < +\infty$, the random variable $Y = g(X)$ has its distribution $p_Y(\mathbf{y})$ satisfying $\mathbb{E}_{p_Y(\mathbf{y})}[\mathbf{y}] = \mathbb{E}_{p_X(\mathbf{x})}[g(\mathbf{x})]$.

Then a generalised version of Jensen's inequality reads as follows.

Proposition 3. (*Generalised Jensen's inequality*) If a function $g(\mathbf{x})$ maps inputs to scalar outputs in \mathbb{R} and $f : \mathbb{R} \rightarrow \mathbb{R}$ is a convex function, then for any distribution $p_X(\mathbf{x})$,

$$\mathbb{E}_{p_X(\mathbf{x})}[f(g(\mathbf{x}))] \geq f(\mathbb{E}_{p_X(\mathbf{x})}[g(\mathbf{x})]),$$

with equality holds iff. f is linear or $p_X(\mathbf{x})$ is a delta measure.

Proof.

$$\begin{aligned} \mathbb{E}_{p_X(\mathbf{x})}[f(g(\mathbf{x}))] &= \mathbb{E}_{p_Y(y)}[f(y)] && \text{(LOTUS applied to } y = g(\mathbf{x})\text{)} \\ &\geq f(\mathbb{E}_{p_Y(y)}[y]) && \text{(Jensen's inequality)} \\ &= f(\mathbb{E}_{p_X(\mathbf{x})}[g(\mathbf{x})]). && \text{(LOTUS applied to } y = g(\mathbf{x})\text{)} \end{aligned}$$

□

Kullback-Leibler (KL) divergence

Kullback-Leibler divergence [Kullback and Leibler, 1951; Kullback, 1959], or *KL divergence*, is arguably one of the most widely used divergence measures in machine learning, statistics, and information theory.

Definition 1. (*Kullback-Leibler Divergence*) The Kullback-Leibler (KL) divergence on \mathcal{P} is defined as a function $\text{KL}[\cdot||\cdot] : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$ with the following form

$$\text{KL}[p||q] = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}, \quad p, q \in \mathcal{P}, \quad (3)$$

where \log is the natural logarithm (to base e).

One can easily check that indeed the above definition is a valid divergence: define $f(x) = -\log x$ (which is convex) and $g(\mathbf{x}) = q(\mathbf{x})/p(\mathbf{x})$, we have

$$\begin{aligned} \text{KL}[p||q] &= \mathbb{E}_{p(\mathbf{x})}[-\log g(\mathbf{x})] \\ &\geq -\log \mathbb{E}_{p(\mathbf{x})}[g(\mathbf{x})] \quad (\text{Jensen's inequality}) \\ &= -\log \int p(\mathbf{x}) \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} = -\log 1 = 0, \end{aligned}$$

and the equality holds iff. $p(\mathbf{x}) = q(\mathbf{x})$.¹ This means one can minimise the KL divergence in order to fit a distribution to a target one. Also notice that the KL divergence is asymmetric, i.e. $\text{KL}[p||q] \neq \text{KL}[q||p]$ in general.

Maximum likelihood estimation (MLE)

Given a dataset $\{(\mathbf{x}_n)\}_{n=1}^N \sim p_{\text{data}}(\mathbf{x})$, we would like to fit to it a generative model $p_{\boldsymbol{\theta}}(\mathbf{x})$ with parameter $\boldsymbol{\theta}$. Since the KL divergence can be used to measure the closeness of the model to the underlying data distribution, it makes sense to find the optimal parameters by minimising the KL divergence:

$$\boldsymbol{\theta}^* = \arg \min \text{KL}[p_{\text{data}}(\mathbf{x})||p_{\boldsymbol{\theta}}(\mathbf{x})]. \quad (4)$$

Expanding the above objective and re-arranging terms, we have

$$\text{KL}[p_{\text{data}}(\mathbf{x})||p_{\boldsymbol{\theta}}(\mathbf{x})] = \underbrace{\mathbb{E}_{p_{\text{data}}(\mathbf{x})}[\log p_{\text{data}}(\mathbf{x})]}_{\text{constant w.r.t. } \boldsymbol{\theta}} - \underbrace{\mathbb{E}_{p_{\text{data}}(\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x})]}_{\text{dependent on } \boldsymbol{\theta}}.$$

This means we can ignore the constant terms w.r.t. $\boldsymbol{\theta}$ and instead work with the following *maximum likelihood* objective:

$$\boldsymbol{\theta}^* = \arg \max \mathbb{E}_{p_{\text{data}}(\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x})]. \quad (5)$$

The obtained optimal parameters $\boldsymbol{\theta}^*$ is called the *maximum likelihood estimate* (MLE) of the parameters. In practice the data distribution is approximated by the empirical distribution on the dataset $\{\mathbf{x}_n\}_{n=1}^N \sim p_{\text{data}}(\mathbf{x})$, leading to

$$\boldsymbol{\theta}^* = \arg \max \frac{1}{N} \sum_{n=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_n). \quad (6)$$

1.2 Variational inference

We are interested in fitting the following latent variable model (LVM) to the data:

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \int p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}. \quad (7)$$

See Figure 1 (a) for a visualisation of the graphical model. In deep generative modelling context, this LVM is often constructed as (for continuous data)

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}), \quad p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; G_{\boldsymbol{\theta}}(\mathbf{z}), \sigma^2 \mathbf{I}), \quad (8)$$

with $G_{\boldsymbol{\theta}}(\cdot)$ define as a neural network transform that is parameterised by weights $\boldsymbol{\theta}$. For discrete variables $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is usually defined as a categorical distribution with a neural network generator in use accordingly. Now to fit $p_{\boldsymbol{\theta}}(\mathbf{x})$ to $p_{\text{data}}(\mathbf{x})$ we optimise the MLE objective (5) w.r.t. $\boldsymbol{\theta}$, which involves computing the integral (7). This is intractable as it involves computing the non-linear transformation $G_{\boldsymbol{\theta}}(\mathbf{z})$ for every single configuration of \mathbf{z} within the support of the Gaussian prior $p(\mathbf{z})$, which is the full space $\mathbf{z} \in \mathbb{R}^d$.

¹Technically speaking: $p(\mathbf{x}) = q(\mathbf{x})$ almost everywhere.

Variational inference provides a variational lower-bound of $\log p_{\theta}(\mathbf{x})$ as an approximation to it. For any distribution $q(\mathbf{z})$ satisfying $q(\mathbf{z}) > 0$ whenever $p_{\theta}(\mathbf{z}|\mathbf{x}) > 0$, we have

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &= \log \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \log \int q(\mathbf{z})\frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z})}d\mathbf{z} \\ &\geq \int q(\mathbf{z})\log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z})}d\mathbf{z} \quad \text{(Jensen's inequality)} \\ &= \mathbb{E}_{q(\mathbf{z})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}[q(\mathbf{z})||p(\mathbf{z})] := \mathcal{L}(\mathbf{x}, q, \theta). \end{aligned} \tag{9}$$

With suitable choice of $q(\mathbf{z})$ and tricks that will be introduced later, this variational lower-bound can be used as a tractable approximation to the marginal log-likelihood $\log p_{\theta}(\mathbf{x})$.

The choice of the $q(\mathbf{z})$ distribution is crucial to the quality of the approximation (or the tightness of the lower-bound). To see this, note that

$$p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_{\theta}(\mathbf{x})}, \quad \text{(Bayes' rule)} \tag{10}$$

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) - \text{KL}[q(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x})] &= \log p_{\theta}(\mathbf{x}) - \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \\ &= \log p_{\theta}(\mathbf{x}) + \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z})p_{\theta}(\mathbf{x})} \right] \quad \text{(Bayes' rule)} \\ &= \mathbb{E}_{q(\mathbf{z})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}[q(\mathbf{z})||p(\mathbf{z})] = \mathcal{L}(\mathbf{x}, q, \theta). \end{aligned} \tag{11}$$

This means the gap (or the approximation error) between the variational lower-bound $\mathcal{L}(\mathbf{x}, q, \theta)$ and the marginal log-likelihood $\log p_{\theta}(\mathbf{x})$ is the KL divergence $\text{KL}[q(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x})]$. Therefore the lower-bound improves as $q(\mathbf{z})$ approaches to the exact posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$. It also motivates the optimisation of the variational lower-bound w.r.t. the q distribution to obtain an approximate posterior: since $\log p_{\theta}(\mathbf{x})$ is constant w.r.t. q , maximising $\mathcal{L}(\mathbf{x}, q, \theta)$ is equivalent to minimising $\text{KL}[q(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x})]$.

1.3 Variational auto-encoders

As discussed so far, we wish to fit the generative model (7) to the data by maximum likelihood (5), and variational inference provides a useful approximation $\mathcal{L}(\mathbf{x}, q, \theta) \leq \log p_{\theta}(\mathbf{x})$ for a given datum \mathbf{x} . Since this approximation is required for every datapoint in $\{\mathbf{x}_n\}_{n=1}^N$, having N separated q distributions $q_1(\mathbf{z}_1), \dots, q_N(\mathbf{z}_N)$ to pair with $\mathbf{x}_1, \dots, \mathbf{x}_N$ can be memory inefficient. However, notice that the exact posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$ depends on the input \mathbf{x} , and the variational lower-bound is tight when $q_n(\mathbf{z}_n) \approx p_{\theta}(\mathbf{z}_n|\mathbf{x}_n)$. This motivates the *variational auto-encoder* (VAE) approach [Kingma and Welling, 2014; Rezende et al., 2014] which defines the q distribution as $q(\mathbf{z}) := q_{\phi}(\mathbf{z}|\mathbf{x})$, with the distribution often defined by a neural network:

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\sigma_{\phi}^2(\mathbf{x}))), \quad \boldsymbol{\mu}_{\phi}(\mathbf{x}), \log \sigma_{\phi}(\mathbf{x}) = \text{NN}_{\phi}(\mathbf{x}). \tag{12}$$

This allows us to define the VAE optimisation objective:

$$\phi^*, \theta^* = \arg \max \mathcal{L}(\phi, \theta), \quad \mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \underbrace{\left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] \right]}_{:= \mathcal{L}(\mathbf{x}, \phi, \theta)}. \tag{13}$$

Analytic KL between factorised Gaussians

Given that both $q_{\phi}(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z})$ are all factorised Gaussian distributions, the KL divergence term in (13) has an analytic form (assuming $\mathbf{z} \in \mathbb{R}^d$):

$$\text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] = \frac{1}{2} (\|\boldsymbol{\mu}_{\phi}(\mathbf{x})\|_2^2 + \|\boldsymbol{\sigma}_{\phi}(\mathbf{x})\|_2^2 - 2\langle \log \boldsymbol{\sigma}_{\phi}(\mathbf{x}), \mathbf{1} \rangle - d). \tag{14}$$

To see this, let us assume two factorised distributions $p(\mathbf{z}) = \prod_{i=1}^d p(z_i)$ and $q(\mathbf{z}) = \prod_{i=1}^d q(z_i)$. Then the KL divergence from q to p can be written as a sum of KL divergences:

$$\begin{aligned} \text{KL}[q(\mathbf{z})||p(\mathbf{z})] &= \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{\prod_{i=1}^d q(z_i)}{\prod_{i=1}^d p(z_i)} \right] = \mathbb{E}_{q(\mathbf{z})} \left[\sum_{i=1}^d \log \frac{q(z_i)}{p(z_i)} \right] \\ &= \sum_{i=1}^d \mathbb{E}_{q(z_i)} \left[\log \frac{q(z_i)}{p(z_i)} \right] = \sum_{i=1}^d \text{KL}[q(z_i)||p(z_i)]. \end{aligned} \quad (15)$$

Then, assuming each $q(z_i)$ and $p(z_i)$ distributions are Gaussians: $q(z_i) = \mathcal{N}(z_i; \mu_i, \sigma_i^2)$, $p(z_i) = \mathcal{N}(z_i; 0, 1)$, we have the KL divergence as:

$$\begin{aligned} \text{KL}[q(z_i)||p(z_i)] &= \mathbb{E}_{q(z_i)} \left[\log \frac{\frac{1}{\sqrt{2\pi\sigma_i^2}} \exp[-\frac{1}{2\sigma_i^2}(z_i - \mu_i)^2]}{\frac{1}{\sqrt{2\pi}} \exp[-\frac{1}{2}z_i^2]} \right] \\ &= \mathbb{E}_{q(z_i)} \left[-\log \sigma_i - \frac{1}{2\sigma_i^2}(z_i - \mu_i)^2 + \frac{1}{2}z_i^2 \right] \\ &= -\log \sigma_i - \frac{1}{2\sigma_i^2} \mathbb{E}_{q(z_i)} [(z_i - \mu_i)^2] + \frac{1}{2} \mathbb{E}_{q(z_i)} [(z_i - \mu_i)^2 - \mu_i^2 + 2\mu_i z_i] \\ &= -\log \sigma_i - \frac{1}{2} + \frac{1}{2}[\sigma_i^2 + \mu_i^2]. \end{aligned} \quad (16)$$

Writing $\boldsymbol{\mu}_\phi(\mathbf{x}) = [\mu_1, \dots, \mu_d]$, $\boldsymbol{\sigma}_\phi(\mathbf{x}) = [\sigma_1, \dots, \sigma_d]$, we can sum up the KL divergence (16) over $i = 1, \dots, d$ and write the resulting $\text{KL}[q(\mathbf{z})||p(\mathbf{z})]$ as (14). This is done by noticing e.g. $\|\boldsymbol{\mu}_\phi(\mathbf{x})\|_2^2 = \sum_{i=1}^d \mu_i^2$ and $\sum_{i=1}^d \log \sigma_i = \langle \log \boldsymbol{\sigma}_\phi(\mathbf{x}), \mathbf{1} \rangle$.

Monte Carlo estimation

The VAE objective $\mathcal{L}(\phi, \theta)$ in (13) is still intractable since the expectation computation $\mathbb{E}_{q_\phi}[\cdot]$ requires evaluating neural network transformations for all possible \mathbf{z} . Monte Carlo (MC) estimation comes into rescue, as we can replace the expectation with MC approximations:

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \approx \log p_\theta(\mathbf{x}|\mathbf{z}), \quad \mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}). \quad (17)$$

By doing so, the gradient of the objective w.r.t. θ can be estimated as

$$\nabla_\theta \mathcal{L}(\mathbf{x}, \phi, \theta) \approx \nabla_\theta \log p_\theta(\mathbf{x}|\mathbf{z}), \quad \mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}). \quad (18)$$

It remains to compute the gradient of the objective w.r.t. ϕ

$$\nabla_\phi \mathcal{L}(\mathbf{x}, \phi, \theta) \approx \nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \nabla_\phi \text{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]. \quad (19)$$

While the gradient w.r.t. the KL term tractable (by differentiate eq. (14) w.r.t. ϕ), MC approximation is still required for the first term in (19).

Reparameterisation trick ²

The MC approximation to $\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$ is further assisted by the reparameterisation trick [Kingma and Welling, 2014; Rezende et al., 2014]. Note that the sampling procedure of a Gaussian variable is the following:

$$\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}) \quad \Leftrightarrow \quad \mathbf{z} = \boldsymbol{\mu}_\phi + \boldsymbol{\sigma}_\phi \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}), \quad (20)$$

²Please note that the reparameterisation trick is not the only method to enable MC estimation of VAE gradients w.r.t. ϕ even when we use Gaussian q distributions. **If interested, see e.g., Section 2.2.3 of this note.*

with \odot denoting element-wise product. Writing $\pi(\epsilon) := \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$ and $T_\phi(\mathbf{x}, \epsilon) := \boldsymbol{\mu}_\phi + \boldsymbol{\sigma}_\phi \odot \epsilon$, we have, by LOTUS,

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{\pi(\epsilon)}[\log p_\theta(\mathbf{x}|T_\phi(\mathbf{x}, \epsilon))], \quad (21)$$

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{\pi(\epsilon)}[\nabla_\phi \log p_\theta(\mathbf{x}|T_\phi(\mathbf{x}, \epsilon))] = \mathbb{E}_{\pi(\epsilon)}[\nabla_\phi \mathbf{z} \nabla_{\mathbf{z}} \log p_\theta(\mathbf{x}|\mathbf{z})|_{\mathbf{z}=T_\phi(\mathbf{x}, \epsilon)}]. \quad (22)$$

Then with MC estimation:

$$\mathbb{E}_{\pi(\epsilon)}[\nabla_\phi \log p_\theta(\mathbf{x}|T_\phi(\mathbf{x}, \epsilon))] \approx \nabla_\phi \mathbf{z} \nabla_{\mathbf{z}} \log p_\theta(\mathbf{x}|\mathbf{z})|_{\mathbf{z}=T_\phi(\mathbf{x}, \epsilon)}, \quad \epsilon \sim \pi(\epsilon). \quad (23)$$

Combined with eq. (18) and mini-batch training, one can compute an MC estimation of the VAE objective (13) as

$$\begin{aligned} \mathcal{L}(\phi, \theta) &\approx \frac{1}{M} \sum_{m=1}^M \log p_\theta(\mathbf{x}_m | T_\phi(\mathbf{x}_m, \epsilon_m)) - \text{KL}[q_\phi(\mathbf{z}_m | \mathbf{x}_m) || p(\mathbf{z}_m)], \\ &\mathbf{x}_1, \dots, \mathbf{x}_m \sim \{\mathbf{x}_n\}^M, \quad \epsilon_1, \dots, \epsilon_M \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \end{aligned} \quad (24)$$

and apply e.g. automatic differentiation to obtain the (MC estimation of) gradient of the VAE objective w.r.t. parameters θ and ϕ .

1.4 Conditional VAE

For conditional generative models, the goal is to generate data (e.g. images) conditioned on additional information. Such additional information can be class labels (which is discrete) or the viewing angle for the image (which is continuous). Mathematically, this corresponds to learning a generative model $p_\theta(\mathbf{x}|\mathbf{y})$ which approximates the data distribution $p_{\text{data}}(\mathbf{x}|\mathbf{y})$. Here \mathbf{x} is the random variable for data (e.g. images) and \mathbf{y} is the random variable corresponding to the additional information (e.g. label or viewing angle).

For the design of the generative model $p_\theta(\mathbf{x}|\mathbf{y})$, we use a conditional LVM as follows:

$$p_\theta(\mathbf{x}|\mathbf{y}) = \int p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y}) p(\mathbf{z}) d\mathbf{z}, \quad (25)$$

See Figure 1 (b) for a visualisation of the graphical model. Often we set $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$. If \mathbf{x} is continuous, then we can define e.g.

$$p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y}) = \mathcal{N}(\mathbf{x}; G_\theta(\mathbf{z}, \mathbf{y}), \sigma^2 \mathbf{I}), \quad (26)$$

with $G_\theta(\mathbf{z}, \mathbf{y})$ defined by a neural network that takes both \mathbf{z} and \mathbf{y} as inputs. Similar to VAEs, learning is done by maximising a variational lower-bound:

$$\phi^*, \theta^* = \arg \max \mathcal{L}(\phi, \theta), \quad \mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{\text{data}}(\mathbf{x}, \mathbf{y})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})} [\log p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y})] - \text{KL}[q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y}) || p(\mathbf{z})]], \quad (27)$$

$$\mathbb{E}_{p_{\text{data}}(\mathbf{x}, \mathbf{y})} [\log p_\theta(\mathbf{x}|\mathbf{y})] \geq \mathcal{L}(\phi, \theta). \quad (28)$$

Although in principle the choice of the q distribution is flexible (since the variational lower-bound holds for almost any q distribution satisfying mild conditions, see Section 1.2), using $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})$ and parameterising it with flexible neural networks would return the best posterior approximation. Using Bayes' rule

$$p_\theta(\mathbf{z}|\mathbf{x}, \mathbf{y}) = \frac{p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y}) p(\mathbf{z})}{p_\theta(\mathbf{x}|\mathbf{y})}, \quad (29)$$

we can show that maximising variational lower-bound w.r.t. q is also equivalent to minimising the KL divergence $\text{KL}[q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})||p_\theta(\mathbf{z}|\mathbf{x}, \mathbf{y})]$:

$$\begin{aligned} & \log p_\theta(\mathbf{x}|\mathbf{y}) - (\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})}[\log p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y})] - \text{KL}[q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})||p(\mathbf{z})]) \\ = & \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{y})q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})}{p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y})p(\mathbf{z})} \right] \\ = & \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})}{p_\theta(\mathbf{z}|\mathbf{x}, \mathbf{y})} \right] = \text{KL}[q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})||p_\theta(\mathbf{z}|\mathbf{x}, \mathbf{y})]. \end{aligned} \quad (30)$$

Therefore if we were to replace $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})$ with $q_\phi(\mathbf{z}|\mathbf{x})$, then the optimal solution does not return the exact posterior approximation, unless the learned generator is degenerate: $G_\theta(\mathbf{z}, \mathbf{y}) = G_\theta(\mathbf{z})$. In such case the \mathbf{y} information is ignored (i.e. $p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{y}) = p_\theta(\mathbf{x}|\mathbf{z})$) and the model is no longer a conditional generative model.

1.5 *Practical interpretations & KL annealing

Comparisons with auto-encoders

Looking at the likelihood part of the VAE objective (13), under Gaussian likelihood assumption we have (by using the reparam. trick)

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{p(\epsilon)} \left[-\frac{1}{2\sigma^2} \|\mathbf{x} - G_\theta(T_\phi(\mathbf{x}, \epsilon))\|_2^2 \right] + \text{const.} \quad (31)$$

On the other hand, an auto-encoder contains a pair of encoder $E_\phi(\cdot)$ and decoder $D_\theta(\cdot)$ which are trained using e.g. ℓ_2 reconstruction loss:

$$\min_{\theta, \phi} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\|\mathbf{x} - D_\theta(E_\phi(\mathbf{x}))\|_2^2]. \quad (32)$$

Comparing the reconstruction loss of the auto-encoder training objective to (31), we see that VAEs can be viewed from a viewpoint of *stochastic* auto-encoder. Architecture-wise, the main difference is the usage of stochastic encoder $T_\phi(\mathbf{x}, \epsilon)$ that injects random noise ϵ to the encoding of \mathbf{x} . Training objective-wise, the VAE objective has the extra $\text{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$ term which regularises the q distribution towards the prior $p(\mathbf{z})$. When $p(\mathbf{z})$ is non-degenerate (e.g. $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$) the resulting $q_\phi(\mathbf{z}|\mathbf{x})$ at optimum is non-degenerate as well, i.e. $\sigma_\phi(\mathbf{x}) > \mathbf{0}$.

KL annealing

Practitioners sometimes find that training VAEs with the original variational lower-bound objective (13) leads to under-fitting issues, in such case often the reconstructed images using the model are blurry. A practical strategy to alleviate this is to introduce a ‘‘KL annealing’’ coefficient β and optimise the θ, ϕ parameters using the following objective:

$$\phi^*, \theta^* = \arg \max \mathcal{L}(\phi, \theta, \beta), \quad \mathcal{L}(\phi, \theta, \beta) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \underbrace{[\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta \text{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]]}_{:=\mathcal{L}(\mathbf{x}, \phi, \theta, \beta)}. \quad (33)$$

If using $0 < \beta < 1$, this objective introduces less regularisation for the $q_\phi(\mathbf{z}|\mathbf{x})$ to be close to the prior $p(\mathbf{z})$. In particular when $\beta = 0$, it results in a stochastic auto-encoder which is trained by the reconstruction loss only. Since stochasticity in \mathbf{z} naturally degrades the quality of reconstruction, training with reconstruction loss only will drive ϕ towards making $\sigma_\phi(\mathbf{x}) \rightarrow \mathbf{0}$ for any \mathbf{x} , which also means $q_\phi(\mathbf{z}|\mathbf{x}) \rightarrow \delta(\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}))$. In such case the resulting model is simply an auto-encoder which cannot be used directly as a generative model for new images.

We should also emphasise that for $\beta < 1$, $\mathcal{L}(\mathbf{x}, \phi, \theta, \beta)$ is no longer a lower-bound for $\log p_\theta(\mathbf{x})$, and the training objective (33) cannot be well justified using (approximate) MLE for learning $p_\theta(\mathbf{x}) \approx$

$p_{\text{data}}(\mathbf{x})$. In fact for small β the learned generative distribution $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ can be very different from $p_{\text{data}}(\mathbf{x})$, again explaining why generation quality can be worse when using such small β values. Therefore β needs to be carefully chosen to achieve the trade of between good reconstruction & good generation. Another strategy is to use different β_t values for different training epochs $t = 0, \dots, T$; a recommended recipe is to select increasing values $0 \leq \beta_1 \leq \dots \leq \beta_T$.

Sometimes $\beta > 1$ values are also used but for a different purpose. Although there is no theoretical guarantee, existing research shows that empirically, with factorised prior $p(\mathbf{z})$ and $\beta > 1$, one can train a VAE to obtain a *disentangled representation*, so that controlled generation can be achieved by varying different dimensions of the \mathbf{z} variable [Higgins et al., 2017].

References

- Bishop, C. M. (2007). *Pattern Recognition and Machine Learning*. Springer.
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. (2017). beta-vae: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*.
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *International Conference on Learning Representations*.
- Kullback, S. (1959). *Information theory and statistics*. John Wiley & Sons.
- Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86.
- Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning*, pages 1278–1286.