

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).

Jensen's inequality

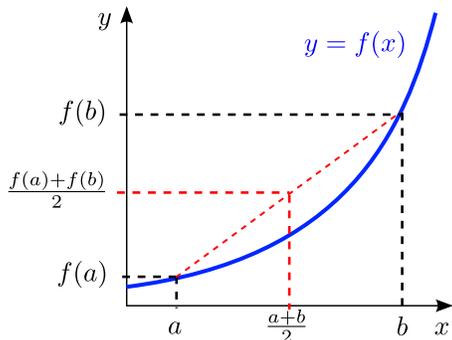
Below we introduce Jensen's inequality as a prerequisite for later discussions on divergences.

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

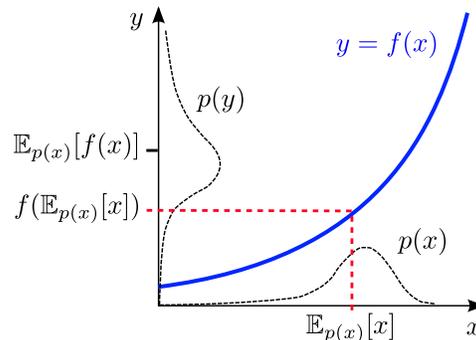
$$\mathbb{E}_{p(x)}[f(x)] \geq f(\mathbb{E}_{p(x)}[x]),$$

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

$$\text{KL}[p||q] = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}, \quad p, q \in \mathcal{P}, \quad (2)$$

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

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

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 $p_{\theta}(\mathbf{x})$ with parameter θ . 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:

$$\theta^* = \arg \min \text{KL}[p_{\text{data}}(\mathbf{x})||p_{\theta}(\mathbf{x})]. \quad (3)$$

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

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

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

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

$$\boldsymbol{\theta}^* = \arg \max \mathbb{E}_{p_{\text{data}}(\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x})]. \quad (4)$$

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

$$\boldsymbol{\theta}^* = \arg \max \frac{1}{N} \sum_{n=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_n). \quad (5)$$

1.2 Variational inference

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

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \int p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}. \quad (6)$$

In deep generative modelling context, this LVM is often constructed as (for continuous data)

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}), \quad p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; G_{\boldsymbol{\theta}}(\mathbf{z}), \sigma^2\mathbf{I}), \quad (7)$$

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 (4) w.r.t. $\boldsymbol{\theta}$, which involves computing the integral (6). 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 $\mathbf{z} \in \mathbb{R}^d$.

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

$$\begin{aligned} \log p_{\boldsymbol{\theta}}(\mathbf{x}) &= \int p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \log \int q(\mathbf{z}) \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z} \\ &\geq \int q(\mathbf{z}) \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z} \quad \text{(Jensen's inequality)} \\ &= \mathbb{E}_{q(\mathbf{z})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \text{KL}[q(\mathbf{z})||p(\mathbf{z})] := \mathcal{L}(\mathbf{x}, q, \boldsymbol{\theta}). \end{aligned} \quad (8)$$

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 $\log p_{\boldsymbol{\theta}}(\mathbf{x})$.

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

$$p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}) = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{x})}, \quad \text{(Bayes' rule)} \quad (9)$$

$$\begin{aligned} \log p_{\boldsymbol{\theta}}(\mathbf{x}) - \text{KL}[q(\mathbf{z})||p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})] &= \log p_{\boldsymbol{\theta}}(\mathbf{x}) - \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{q(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})} \right] \\ &= \log p_{\boldsymbol{\theta}}(\mathbf{x}) + \mathbb{E}_{q(\mathbf{z})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{x})} \right] \quad \text{(Bayes' rule)} \\ &= \mathbb{E}_{q(\mathbf{z})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \text{KL}[q(\mathbf{z})||p(\mathbf{z})] = \mathcal{L}(\mathbf{x}, q, \boldsymbol{\theta}). \end{aligned} \quad (10)$$

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

1.3 Variational auto-encoders

As discussed so far, we wish to fit the generative model (6) to the data by maximum likelihood (4), and variational inference provides a useful approximation $\mathcal{L}(\mathbf{x}, q, \boldsymbol{\theta}) \leq \log p_{\boldsymbol{\theta}}(\mathbf{x})$ 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 $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ depends on the input \mathbf{x} , and the variational lower-bound is tight when $q_n(\mathbf{z}_n) \approx p_{\boldsymbol{\theta}}(\mathbf{z}_n|\mathbf{x}_n)$. 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:

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\sigma_{\phi}^2(\mathbf{x}))), \quad \boldsymbol{\mu}_{\phi}(\mathbf{x}), \log \sigma_{\phi}(\mathbf{x}) = \text{NN}_{\phi}(\mathbf{x}). \quad (11)$$

This allows us to define the VAE optimisation objective:

$$\phi^*, \boldsymbol{\theta}^* = \arg \max \mathcal{L}(\phi, \boldsymbol{\theta}), \quad \mathcal{L}(\phi, \boldsymbol{\theta}) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \underbrace{[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]]}_{:= \mathcal{L}(\mathbf{x}, \phi, \boldsymbol{\theta})}. \quad (12)$$

Given that both $q_{\phi}(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z})$ are all factorised Gaussian distributions, the KL divergence term in (12) has an analytic form (assuming $\mathbf{z} \in \mathbb{R}^d$):

$$\text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] = \frac{1}{2} (\|\boldsymbol{\mu}_{\phi}(\mathbf{x})\|_2^2 + \|\boldsymbol{\sigma}_{\phi}(\mathbf{x})\|_2^2 - 2\langle \log \sigma_{\phi}(\mathbf{x}), \mathbf{1} \rangle - d). \quad (13)$$

Monte Carlo estimation

The VAE objective $\mathcal{L}(\phi, \boldsymbol{\theta})$ in (12) 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:

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] \approx \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}), \quad \mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}). \quad (14)$$

By doing so, the gradient of the objective w.r.t. $\boldsymbol{\theta}$ can be estimated as

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}, \phi, \boldsymbol{\theta}) \approx \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}), \quad \mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}). \quad (15)$$

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

$$\nabla_{\phi} \mathcal{L}(\mathbf{x}, \phi, \boldsymbol{\theta}) \approx \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \nabla_{\phi} \text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]. \quad (16)$$

While the gradient w.r.t. the KL term tractable (by differentiate eq. (13) w.r.t. ϕ), MC approximation is still required for the first term in (16).

Reparameterisation trick

The MC approximation to $\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]$ 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:

$$\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}) \quad \Leftrightarrow \quad \mathbf{z} = \boldsymbol{\mu}_{\phi} + \boldsymbol{\sigma}_{\phi} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}), \quad (17)$$

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

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{\pi(\boldsymbol{\epsilon})}[\log p_\theta(\mathbf{x}|T_\phi(\mathbf{x}, \boldsymbol{\epsilon}))], \quad (18)$$

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{\pi(\boldsymbol{\epsilon})}[\nabla_\phi \log p_\theta(\mathbf{x}|T_\phi(\mathbf{x}, \boldsymbol{\epsilon}))] = \mathbb{E}_{\pi(\boldsymbol{\epsilon})}[\nabla_\phi \mathbf{z} \nabla_{\mathbf{z}} \log p_\theta(\mathbf{x}|\mathbf{z})|_{\mathbf{z}=T_\phi(\mathbf{x}, \boldsymbol{\epsilon})}]. \quad (19)$$

Then with MC estimation:

$$\mathbb{E}_{\pi(\boldsymbol{\epsilon})}[\nabla_\phi \log p_\theta(\mathbf{x}|T_\phi(\mathbf{x}, \boldsymbol{\epsilon}))] \approx \nabla_\phi \mathbf{z} \nabla_{\mathbf{z}} \log p_\theta(\mathbf{x}|\mathbf{z})|_{\mathbf{z}=T_\phi(\mathbf{x}, \boldsymbol{\epsilon})}, \quad \boldsymbol{\epsilon} \sim \pi(\boldsymbol{\epsilon}) \quad (20)$$

Combined with eq. (15) and mini-batch training, one can compute an MC estimation of the VAE objective (12) as

$$\begin{aligned} \mathcal{L}(\phi, \theta) &\approx \frac{1}{M} \sum_{m=1}^M \log p_\theta(\mathbf{x}_m | T_\phi(\mathbf{x}_m, \boldsymbol{\epsilon}_m)) - \text{KL}[q_\phi(\mathbf{z}_m | \mathbf{x}_m) || p(\mathbf{z}_m)], \\ &\quad \mathbf{x}_1, \dots, \mathbf{x}_m \sim \{\mathbf{x}_n\}^M, \quad \boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_M \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \end{aligned} \quad (21)$$

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

References

- 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.