

Basics VAE

April 21, 2023

Useful papers - [\[Doe16\]](#); [\[KW13\]](#)

Contents

1	EM-algorithm	2
1.1	Problem statement	2
1.2	Gaussian mixture. Likelihoods comparison	2
1.3	Main algorithm	3
1.4	Gaussian mixture. Complete solution	3
2	Introduction	4
2.1	Intuition behind VAE	4
2.2	ELBO derivation	5
2.3	Some algebra	8
2.4	Final optimization problem	10
2.5	Prior hole problem	12
3	Connection between VAE and EM-algorithm	12
4	Information theory	13

1 EM-algorithm

1.1 Problem statement

As usually, we have a parameterized model – $p_\theta(x)$. Our goal is to find parameters θ in such a way $p_\theta(x) \approx p_{\text{data}}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x_i)$, where $\mathbf{X} = \{x_i\}_{i=1}^N$ - our data. We can use MLE for this purpose:

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} \mathbb{E}_{p_{\text{data}}(x)} \log p_\theta(x) = \operatorname{argmax}_{\theta} \mathcal{L}(\theta|\mathbf{X}) \\ \mathcal{L}(\theta|\mathbf{X}) &= \mathbb{E}_{p_{\text{data}}(x)} \log p_\theta(x) = \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \log p_\theta(x)\end{aligned}\tag{1}$$

However, in most cases $\mathcal{L}(\theta|\mathbf{X})$ is hard to calculate. It is stated that the following likelihood is simpler to calculate:

$$\mathcal{L}(\theta|\mathbf{X}, \mathbf{Z}) = \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z)} \log p_\theta(x, z)\tag{2}$$

Here we added latent variables, z . In this sense $\mathcal{L}(\theta|\mathbf{X})$ calls as a incomplete likelihood and $\mathcal{L}(\theta|\mathbf{X}, \mathbf{Z})$ as complete. However, it is not clear that the complete likelihood is simpler to calculate than incomplete. Let's consider the following example.

1.2 Gaussian mixture. Likelihoods comparison

Let we need to approximate the mixture of the gaussians. Then

$$\begin{aligned}p_\theta(x) &= \sum_{j=1}^k w_j \mathcal{N}(x|\alpha_j), \alpha_j = (\mu_j, \sigma_j^2) \\ \theta &= (w_1, \dots, w_k, \alpha_1, \dots, \alpha_k)\end{aligned}\tag{3}$$

Let's compare two likelihoods

$$\mathcal{L}(\theta|\mathbf{X}) = \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \log \sum_{j=1}^k w_j \mathcal{N}(x|\alpha_j)\tag{4}$$

We can define latent variables as a components of the gaussians, i.e. $z_i \in \{1, \dots, k\}$.

$$\begin{aligned}\mathcal{L}(\theta|\mathbf{X}, \mathbf{Z}) &= \mathbb{E}_{p_{\text{data}}(x)} \mathbb{E}_{q(z)} \log p_\theta(x, z) \\ &= \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \mathbb{E}_{q(z)} [\log p_\theta(x|z) + \log q(z)] \\ &= \max_{\theta} \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \mathbb{E}_{q(z)} \log p_\theta(x|z)\end{aligned}\tag{5}$$

To calculate $\mathcal{L}(\theta|\mathbf{X}, \mathbf{Z})$ we need to know $\log p_\theta(x|z)$ and $\log q(z)$. However, the first term is much simpler than $\log p_\theta(x)$. Indeed, $\log p_\theta(x|j) = \log \mathcal{N}(x|\alpha_j)$ insted of sum, which was in $p_\theta(x)$. As we further see, EM-algorithm will provide the access to $q(z)$ in

the closed-form. Thus, we have that the complete likelihood is easier to calculate than incomplete

$$\begin{aligned}\mathcal{L}(\theta|\mathbf{X}) &= \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \log \sum_{j=1}^k w_j \mathcal{N}(x|\alpha_j) \\ \mathcal{L}(\theta|\mathbf{X}, \mathbf{Z}) &= \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \mathbb{E}_{q(j)} \log \mathcal{N}(x|\alpha_j)\end{aligned}\tag{6}$$

1.3 Main algorithm

So, we stated that we want to work with $\mathcal{L}(\theta|\mathbf{X}, \mathbf{Z})$ instead of $\mathcal{L}(\theta|\mathbf{X})$. But how to connect them? In can be shown that

$$\log p_\theta(x) = \mathbb{E}_{q(z)} \log p(x, z|\theta) - \mathbb{E}_{q(z)} \log q(z) + \mathbb{KL}(q(z)||p_\theta(z|x))\tag{7}$$

We know that $\mathbb{KL}(q(z)||p_\theta(z|x)) \geq 0$. This fact immediately gives us the form of the $q(z)$. In other words, it should equal to the true posterior, i.e. $p_\theta(z|x)$ (we will use Bayes theorem to calculate it). Thus, EM-algorithm is the following (Expectation-Maximization algorithm)

E-step:

$$q(z) = p_{\theta_{old}}(z|x) = \frac{p_{\theta_{old}}(x, z)}{\int p_{\theta_{old}}(x, z) dz}\tag{8}$$

M-step:

$$\theta_{new} = \operatorname{argmax}_{\theta} \mathbb{E}_{p_{\mathbf{data}}(x)} \mathbb{E}_{q(z)} \log p_\theta(x, z)\tag{9}$$

Two questions may arise:

- 1) Why E-step is E-step? There is no expectation.
- 2) What is it θ_{old} ?

Originally, at the E-step we are interested in the calculation of the expectation $\mathbb{E}_{q(z)} \log p_\theta(x, z)$. However, since the main problem here is to find $q(z)$, we focus only on $q(z)$. Then we can just put them together as $\mathbb{E}_{q(z)} \log p_\theta(x, z) = \int p_{\theta_{old}}(z|x) \log p_\theta(x, z) dz$. θ_{old} is parameters from the previous step.

It is necessary to highlight one specific point. At the E-step we have to calculate normalizing constant, $\int p_{\theta_{old}}(x, z) dz = p_{\theta_{old}}(x)$. As we can see, it is a part of the initial incomplete likelihood, but previously we said that in most cases it is hard to calculate it. Actually, in practice we do not need to calculate $q(z)$ in the closed form, since eq. (9) requires only expectation by this distribution. Usually, we will need only statistics of the $q(z)$, e.g. $\mathbb{E}_{q(z)} z$ or $\mathbb{E}_{q(z)} z^2$.

1.4 Gaussian mixture. Complete solution

Lets back to the

$$\begin{aligned}p_\theta(x) &= \sum_{j=1}^k w_j \mathcal{N}(x|\alpha_j), \alpha_j = (\mu_j, \sigma_j^2) \\ \theta &= (w_1, \dots, w_k, \alpha_1, \dots, \alpha_k), z_k \in \{1, \dots, k\} = \mathbf{Z}\end{aligned}\tag{10}$$

So, E-step is the following:

$$q(z_i) = \frac{p_{\theta_{old}}(x, z_i)}{\int p_{\theta_{old}}(x, z) dz} = \frac{p_{\theta_{old}}(x|z_i)p(z_i)}{p_{\theta_{old}}(x)} = \frac{\mathcal{N}(x|\alpha_i^{old})w_i^{old}}{\sum_{j=1}^k w_j^{old}\mathcal{N}(x|\alpha_j^{old})} \quad (11)$$

Then we have to make M-step:

$$\begin{aligned} \theta_{new} &= \operatorname{argmax}_{\theta} \mathbb{E}_{p_{\mathbf{data}}(x)} \mathbb{E}_{q(z)} \log p_{\theta}(x, z) \\ &= \operatorname{argmax}_{\theta} \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \sum_{z_i \in \mathbf{Z}} q(z_i) [\log w_i + \log \mathcal{N}(x|\alpha_i)] \end{aligned} \quad (12)$$

Note that maximization with respect to the w_i is a conditional optimization problem since $\sum_i w_i = 1$. Thus we have to optimize the following Lagrangian:

$$\begin{aligned} \nabla_{w_j} \left[\frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \sum_{z_i \in \mathbf{Z}} q(z_i) [\log w_i + \log \mathcal{N}(x|\alpha_i)] - \lambda \left(\sum_k w_k - 1 \right) \right] &= 0 \\ \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} q(z_j) \frac{1}{w_j} - \lambda &= 0 \\ w_i^{new} &= \frac{1}{\lambda} \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} q(z_j), \\ \lambda &= 1 \text{ (can be shown using normalization)} \end{aligned} \quad (13)$$

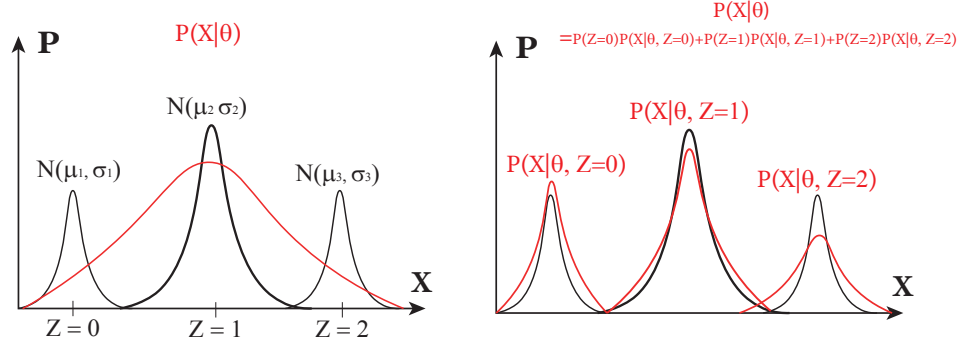
Then we have to find α_i . It can be found in the same spirit.

$$\nabla_{\alpha_j} \left[\frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \sum_{z_i \in \mathbf{Z}} q(z_i) [\log w_i + \log \mathcal{N}(x|\alpha_i)] \right] = 0 \quad (14)$$

2 Introduction

2.1 Intuition behind VAE

Let x_1, \dots, x_n - independent and identically distributed (i.i.d) random variables from distribution $P_{data}(X)$ (let's say this would be training dataset). Goal of any generative model is to reproduce the distribution $P_{data}(X)$ (or at least be able to make samples from it). VAE works with latent variables Z from some distribution $P(Z)$. The main idea behind latent distribution is to reproduce initial object (X variable) using hidden description (Z). Note that we also assume that variable Z can be easily sampled from $P(Z)$. Let's discuss the motivation of this idea in more detail. It is quite logical to assume the following fact: if we have some additional information about variable X , then the conditional distribution $P_{data}(X|Z)$ is simpler than initial $P_{data}(X)$. Let's show this with an example of mixture of gaussians (Figure 2.1). In this example we have mix of three gaussians. If we will approximate it directly using normal distribution as $P(X|\Theta)$, then we obtain bad approximation (as can be seen from figure). However, if we add latent description Z that define each gaussian, we can consider conditional distribution as gaussian. After that, using marginalization, we obtain $P(X|\Theta)$ with higher approximation accuracy to $P_{data}(X)$. Thus, conditional distribution much simpler than $P(X|\Theta)$ and final approximation quite better than non-latent approach.



2.2 ELBO derivation

So, if we have latent description, we need to build parameterized function $f : Z \times \Theta \rightarrow X$ to display hidden variable to the initial. In this way we create samples from conditional distribution $P(X|Z, \Theta)$. We can obtain marginal distribution as

$$\begin{aligned}
 P(X|\Theta) &= \int_Z P(X, Z|\Theta) dz = \\
 &= [P(X, Z|\Theta) = P(Z, X|\Theta) = P(X|Z, \Theta)P(Z|\Theta), P(Z|\Theta) = P(Z)] = \quad (15) \\
 &= \int_Z P(X|Z, \Theta)P(Z)dz = \mathbb{E}_{z \sim P} P(X|Z, \Theta) \approx \frac{1}{m} \sum_{i=1}^m P(X|z_i, \Theta)
 \end{aligned}$$

Here, we used Monte-Carlo approximation of the integral, thus $z_i \sim P(Z)$. This marginal will be an approximation of our data distribution $P_{data}(X)$. VAE using following assumptions: $P(X|Z, \Theta) = \mathbf{N}(X|f_{\Theta}(Z), \sigma I)$, $P(Z) = \mathbf{N}(Z|0, I)$. Usually, if we have some function that want to approximate, lets use universal neural network for this purpose. That is, $f_{\Theta}(Z)$ is neural network. We can use maximum marginal likelihood estimation for approximation the data distribution:

$$\begin{aligned}
 \Theta^* &= \operatorname{argmax}_{\Theta} P(X|\Theta) = \operatorname{argmax}_{\Theta} \mathbb{E}_{z \sim P} P(X|Z, \Theta) \\
 &= \operatorname{argmax}_{\Theta} \int dZ P(Z)P(X|Z, \Theta) = \operatorname{argmax}_{\Theta} \prod_j \int dz_j P(z_j)P(x_j|z_j, \Theta) \\
 &\approx \operatorname{argmax}_{\Theta} \frac{1}{m} \sum_{i=1}^m \prod_{j=1}^n P(x_j|z_{ji}, \Theta) \quad (16) \\
 P_{data}(X) &\approx P(X|\Theta^*)
 \end{aligned}$$

A question may arise: why this optimization problem gives the approximation of the $P_{data}(X)$. Lets show it. From the properties of the KL-divergence (in Section 4 details are provided): if $\mathbb{KL}(P_{data}(X), P(X|\Theta)) = 0$, then $P(X|\Theta) = P_{data}(X)$. We need to show that minimization of the KL-divergence is equivalence to maximum likelihood estimation.

$$\begin{aligned}
\operatorname{argmin}_{\Theta} \mathbb{KL}(P_{data}(X), P(X|\Theta)) &= \operatorname{argmin}_{\Theta} \int dX P_{data}(X) \log \frac{P_{data}(X)}{P(X|\Theta)} = \\
\operatorname{argmin}_{\Theta} \left[\int dX P_{data}(X) \log P_{data}(X) - \int dX P_{data}(X) \log P(X|\Theta) \right] &= \\
= \operatorname{argmax}_{\Theta} \int dX P_{data}(X) \log P(X|\Theta) &= \operatorname{argmax}_{\Theta} \mathbb{E}_{X \sim P_{data}} \log P(X|\Theta)
\end{aligned} \tag{17}$$

So, we show it for logarithm. But the maximum of $\log P(X|\Theta)$ is equal to maximum of $P(X|\Theta)$ due to monotonic properties of the logarithm. However, we also have expectation on samples from data distribution $\mathbb{E}_{X \sim P_{data}}$. This expectation is sometimes called an experimental distribution, and it is discrete, not continuous, as I wrote in (17) (but im not sure :)). Then,

$$\begin{aligned}
\operatorname{argmax}_{\Theta} \mathbb{E}_{X \sim P_{data}} \log P(X|\Theta) &= \operatorname{argmax}_{\Theta} \sum_{x_i} P_{data}(X = x_i) \log P(X = x_i|\Theta) = \\
&= \operatorname{argmax}_{\Theta} \frac{1}{n} \sum_i \log P(x_i|\Theta) = \operatorname{argmax}_{\Theta} \log P(X|\Theta)
\end{aligned} \tag{18}$$

Here, x_i - sample from dataset. We used that experimental distribution is uniform over the number of dataset samples, i.e. $i \sim \mathbb{U}(1, \dots, n)$. Hence, we have shown the equivalence between KL-divergence and log-likelihood maximization. Then, we can approximate data distribution using (16).

So, as was said, we can use log-likelihood maximization. But in (16) we have a product of distribution without logarithm. It is so hard to take gradients from product, this procedure is ineffective and resource-consuming. To avoid a product lets change it to logarithm

$$\Theta^* = \operatorname{argmax}_{\Theta} \log P(X|\Theta) = \operatorname{argmax}_{\Theta} \log \mathbb{E}_{z \sim P} P(X|Z, \Theta) \tag{19}$$

Unfortunately, we have met a problem. We cant forward logarithm through expectation ($\log \mathbb{E} \neq \mathbb{E} \log$), thus cant avoid a product. For clarity lets prove $\log \mathbb{E} \neq \mathbb{E} \log$. As we know, logarithm is concave function. Function f calls concave if

$$f(y) \leq f(x) + f'(x)(y - x), \quad \forall y \tag{20}$$

It means that the tangent at any point lies above the function, x is a random but fixed point. Let y and x be a random variables. If we can choose x the way we want, let $x = \mathbb{E} y$. Taking expectations of both parts:

$$\begin{aligned}
\mathbb{E} f(y) &\leq f(\mathbb{E} x) + \mathbb{E} f'(x)(y - x) \\
\mathbb{E} f(y) &\leq f(\mathbb{E} y) + f'(\mathbb{E} y)(\mathbb{E} y - \mathbb{E} y) = f(\mathbb{E} y)
\end{aligned} \tag{21}$$

It calls Jensen inequality, lets do it for (19)

$$\operatorname{argmax}_{\Theta} \log \mathbb{E}_{z \sim P} P(X|Z, \Theta) \geq \operatorname{argmax}_{\Theta} \mathbb{E}_{z \sim P} \log P(X|Z, \Theta) \tag{22}$$

Finally, we did everything and ready for computational optimization!... No. So, we can maximize $\mathbb{E}_{z \sim P} \log P(X|Z, \Theta)$, but the gap between $\mathbb{E}_{z \sim P} \log P(X|Z, \Theta)$ and $\log \mathbb{E}_{z \sim P} P(X|Z, \Theta)$ can be very large, thus we obtain bad parameter Θ . We have to consider another approach. More precisely, we need something that decrease mentioned gap. This is where VAE comes in.

VAE introduce new distribution - $Q(Z)$. As we will further, this $Q(Z)$ will help to decrease the gap. But, now we need to connect $Q(Z)$ and $P(X|\Theta)$.

$$\log P(X|\Theta) = \int_Z Q(Z) \log P(X|\Theta) dZ = [P(X|\Theta) = \frac{P(X, Z|\Theta)}{P(Z|X, \Theta)}] = \quad (23)$$

$$= \int_Z Q(Z) \log \frac{P(X, Z|\Theta)}{P(Z|X, \Theta)} dZ = \int_Z Q(Z) \log \frac{P(X, Z|\Theta)Q(Z)}{P(Z|X, \Theta)Q(Z)} dZ$$

$$= \int_Z Q(Z) \log \frac{P(X, Z|\Theta)}{Q(Z)} dZ + \int_Z Q(Z) \log \frac{Q(Z)}{P(Z|X, \Theta)} dZ = \quad (24)$$

$$= \int_Z Q(Z) \log P(X|Z, \Theta) dZ + \int_Z Q(Z) \log \frac{P(Z)}{Q(Z)} dZ + \int_Z Q(Z) \log \frac{Q(Z)}{P(Z|X, \Theta)} dZ$$

Finally, we have

$$\log P(X|\Theta) - \mathbb{KL}(Q(Z)||P(Z|X)) = \mathbb{E}_{z \sim Q} \log P(X|Z, \Theta) - \mathbb{KL}(Q(Z)||P(Z)) = \mathcal{L}(\Theta, Q)$$

$$\log P(X|\Theta) = \mathcal{L}(\Theta, Q) + \mathbb{KL}(Q(Z)||P(Z|X, \Theta)) \quad (25)$$

$\mathcal{L}(\Theta, Q)$ calls evidence lower bound (ELBO).

Definition 1 *Variational lower bound*

Function $g(x, y(x))$ calls lower bound for function $f(x)$ if and only if

1. $\forall x, y \quad g(x, y(x)) \leq f(x)$
2. $\exists x_0 : g(x_0, y(x_0)) = f(x_0)$

As can be seen $\mathcal{L}(\Theta, Q)$ satisfy conditionals of definition:

$$\begin{aligned} &1. \log P(X|\Theta) \geq \mathcal{L}(\Theta, Q); \\ &\mathbb{KL}(Q(Z)||P(Z|X)) \geq 0, \forall \Theta, \Phi \end{aligned} \quad (26)$$

$$2. \exists Q : Q(Z) = P(Z|X, \Theta) \Rightarrow \log P(X|\Theta) = \mathcal{L}(\Theta, Q)$$

Lets look at (25) more precisely. We found "something" that can decrease the gap, it is Q . Indeed, if we maximize $\mathcal{L}(\Theta, Q)$ with respect to Q then $\mathbb{KL}(Q(Z)||P(Z|X, \Theta))$ should decrease, because $\log P(X|\Theta)$ does not depend on Q . Thus, maximization of $\mathcal{L}(\Theta, Q)$ with respect to Q decrease the gap between $\log P(X|\Theta)$ and $\mathcal{L}(\Theta, Q)$. Also, maximization by Θ moves $P(X|\Theta)$ towards $P_{data}(X)$. So, ELBO optimization is a very

nice one!

$$\log P(X|\Theta) = \sum_{i=1}^n \log p(x_i, \Theta) \geq \sum_{i=1}^n \left[\int_{z_i} Q(z_i) \log p(x_i|z_i, \Theta) dz_i - \mathbb{KL}(Q(z_i)||p(z_i)) \right] \quad (27)$$

Where n is the number of training data. So, the main idea of VAE is to maximize ELBO instead of maximizing likelihood directly, because $\mathbb{KL}(Q(Z)||P(Z|X, \Theta))$ is intractable, because we don't know $P(Z|X, \Theta)$. Also, it is hard to optimize by function. Thus, let's constrain Q by parameterical class $Q(z_i|\phi_i) = \mathbf{N}(z_i|\mu(\phi_i), \sigma(\phi_i))$ and consider $P(z_i) = \mathbf{N}(z_i|0, I)$. Here ϕ_i depends on the number of object, because for every object we want to get distribution of latent variable. We have two more problems here: 1) Using stochastic gradient descent, we will not update all parameters at the end of the epoch, thus it leads to slow convergence; 2) If we want to calculate latent distribution for a new object, we have to retrain our network, because we have no parameters ϕ_i for a new object. To solve both problems we can make dependencies on X and no dependencies on i , i.e. $Q(z_i|x_i, \phi)$. So, now ELBO looks like

$$\mathcal{L}(\Theta, \phi) = \sum_{i=1}^n \left[\int_{z_i} Q(z_i|x_i, \phi) \log p(x_i|z_i, \Theta) dz_i - \mathbb{KL}(Q(z_i|x_i, \phi)||p(z_i)) \right] \quad (28)$$

With assumptions:

$$\begin{aligned} Q(z_i|x_i, \phi) &= \mathbf{N}(z_i|\mu(\phi, x_i), \sigma(\phi, x_i)) - \text{Encoder} \\ p(z_i) &= \mathbf{N}(z_i|0, I) - \text{Prior latent distribution} \\ p(x_i|z_i, \Theta) &= \mathbf{N}(x_i|f(z_i, \Theta), s^2 I) - \text{Decoder} \\ x_i, f(z_i, \Theta) &\in \mathbb{R}^D; \mu, \sigma, z_i \in \mathbb{R}^d \\ \phi, \Theta &- \text{Parameters of neural network} \\ (x_1, \dots, x_n) &- \text{Training data} \end{aligned} \quad (29)$$

2.3 Some algebra

Let's calculate the second term of ELBO - $\mathbb{KL}(Q(z_i|x_i, \phi)||P(z_i))$ (will be tough)

$$\begin{aligned} \mathbb{KL}(\mathbf{N}(z_i|\mu(x_i, \phi), \sigma(x_i, \phi))||\mathbf{N}(z_i|0, I)) &= \int_{z_i} \mathbf{N}(z_i|\mu, \sigma) \log \frac{\mathbf{N}(z_i|\mu, \sigma)}{\mathbf{N}(z_i|0, I)} dz_i = \\ &= \int_{z_i} \mathbf{N}(z_i|\mu, \sigma) \log \mathbf{N}(z_i|\mu, \sigma) dz_i - \int_{z_i} \mathbf{N}(z_i|\mu, \sigma) \log \mathbf{N}(z_i|0, I) dz_i = \mathbf{I}_1 - \mathbf{I}_2 \end{aligned} \quad (30)$$

Note that $\mathbf{cov}(z_i) = \text{diag}(\sigma)$

$$\begin{aligned}
\mathbf{N}(z_i|\mu, \sigma) &= \frac{1}{(2\pi)^{\frac{d}{2}} |\mathbf{cov}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (z_i - \mu)^T \mathbf{cov}^{-1} (z_i - \mu) \right) = \\
&= \frac{1}{(2\pi)^{\frac{d}{2}} \prod_{k=1}^d \sigma_k} \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j - \mu_j)^2 \frac{1}{\sigma_j^2} \right) \\
\mathbf{N}(z_i|0, I) &= \frac{1}{(2\pi)^{\frac{d}{2}}} \exp \left(-\frac{1}{2} z_i^T z_i \right) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j)^2 \right)
\end{aligned} \tag{31}$$

Let's consider \mathbf{I}_2

$$\begin{aligned}
\mathbf{I}_2 &= \frac{1}{(2\pi)^{\frac{d}{2}} \prod_{k=1}^d \sigma_k} \cdot \int_{z_i} \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j - \mu_j)^2 \frac{1}{\sigma_j^2} \right) \log \left(\frac{1}{(2\pi)^{\frac{d}{2}}} \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j)^2 \right) \right) dz_i = \frac{1}{(2\pi)^{\frac{d}{2}} \prod_{k=1}^d \sigma_k} \cdot \\
&\cdot \left[\int_{z_i} \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j - \mu_j)^2 \frac{1}{\sigma_j^2} \right) \log \left(\frac{1}{(2\pi)^{\frac{d}{2}}} \right) dz_i - \frac{1}{2} \int_{z_i} \sum_{j=1}^d (z_i^j)^2 \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j - \mu_j)^2 \frac{1}{\sigma_j^2} \right) dz_i \right] = \\
&= \log \left(\frac{1}{(2\pi)^{\frac{d}{2}}} \right) - \frac{1}{2(2\pi)^{\frac{d}{2}} \prod_{k=1}^d \sigma_k} \mathbf{J} \\
&\tag{32} \\
\mathbf{J} &= \int_{z_i} \sum_{j=1}^d (z_i^j)^2 \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j - \mu_j)^2 \frac{1}{\sigma_j^2} \right) dz_i = \int_{z_i} (z_i^1)^2 \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j - \mu_j)^2 \frac{1}{\sigma_j^2} \right) dz_i + \dots + \\
&+ \int_{z_i} (z_i^d)^2 \exp \left(-\frac{1}{2} \sum_{j=1}^d (z_i^j - \mu_j)^2 \frac{1}{\sigma_j^2} \right) dz_i = \mathbf{J}_1 + \dots + \mathbf{J}_n \\
&\tag{33}
\end{aligned}$$

$$\begin{aligned}
\mathbf{J}_1 &= \int_{-\infty}^{\infty} (z_i^1)^2 \exp\left(-\frac{1}{2}(z_i^1 - \mu_1)^2 \frac{1}{\sigma_1^2}\right) dz_i^1 \dots \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2}(z_i^d - \mu_d)^2 \frac{1}{\sigma_d^2}\right) dz_i^d = \\
&= \int_{-\infty}^{\infty} (z_i^1)^2 \exp\left(-\frac{1}{2}(z_i^1 - \mu_1)^2 \frac{1}{\sigma_1^2}\right) dz_i^1 (2\pi)^{\frac{d-1}{2}} \prod_{k=2}^d \sigma_k = \\
&= (2\pi)^{\frac{d-1}{2}} \prod_{k=2}^d \sigma_k \int_{-\infty}^{\infty} ((z_i^1 - \mu_1)^2 + 2\mu_1 z_i^1 - \mu_1^2) \exp\left(-\frac{1}{2}(z_i^1 - \mu_1)^2 \frac{1}{\sigma_1^2}\right) dz_i^1 = \\
&= (2\pi)^{\frac{d-1}{2}} \prod_{k=2}^d \sigma_k \left[\int_{-\infty}^{\infty} (z_i^1 - \mu_1)^2 \exp\left(-\frac{1}{2}(z_i^1 - \mu_1)^2 \frac{1}{\sigma_1^2}\right) dz_i^1 + 2\mu_1 \int_{-\infty}^{\infty} z_i^1 \exp\left(-\frac{1}{2}(z_i^1 - \mu_1)^2 \frac{1}{\sigma_1^2}\right) dz_i^1 \right] - \\
&\quad - (2\pi)^{\frac{d-1}{2}} \prod_{k=2}^d \sigma_k \mu_1^2 \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2}(z_i^1 - \mu_1)^2 \frac{1}{\sigma_1^2}\right) dz_i^1 = (2\pi)^{\frac{d}{2}} \prod_{k=1}^d \sigma_k (\sigma_1^2 + \mu_1^2)
\end{aligned} \tag{34}$$

$$\mathbf{J} = (2\pi)^{\frac{d}{2}} \prod_{k=1}^d \sigma_k \sum_{j=1}^d (\sigma_j^2 + \mu_j^2) \tag{35}$$

$$\mathbf{I}_2 = \log\left(\frac{1}{(2\pi)^{\frac{d}{2}}}\right) - \frac{1}{2} \sum_{j=1}^d (\sigma_j^2 + \mu_j^2)$$

Analogue

$$\mathbf{I}_1 = \log\left(\frac{1}{(2\pi)^{\frac{d}{2}}}\right) - \frac{1}{2} \sum_{j=1}^d (1 + \log \sigma_j^2) \tag{36}$$

Finally

$$\begin{aligned}
&\mathbb{KL}(\mathbf{N}(z_i | \mu(x_i, \phi), \sigma(x_i, \phi)) || \mathbf{N}(z_i | 0, I)) = \\
&= -\frac{1}{2} \sum_{j=1}^d (1 + \log \sigma_j(x_i, \phi)^2 - \sigma_j(x_i, \phi)^2 - \mu_j(x_i, \phi)^2)
\end{aligned} \tag{37}$$

2.4 Final optimization problem

Our goal is to find a gradient for maximizing ELBO

$$\begin{aligned}
&\nabla_{\Theta, \phi} \mathcal{L}(\Theta, \phi) = \\
&= \nabla_{\Theta, \phi} \sum_{i=1}^n \left[\int_{z_i} Q(z_i | x_i, \phi) \log p(x_i | z_i, \Theta) dz_i + \frac{1}{2} \sum_{j=1}^d (1 + \log \sigma_j(x_i, \phi)^2 - \sigma_j(x_i, \phi)^2 - \mu_j(x_i, \phi)^2) \right] = \\
&= \sum_{i=1}^n \nabla_{\Theta, \phi} \mathcal{L}_i(\Theta, \phi) - ?
\end{aligned} \tag{38}$$

We will use stochastic gradient for gradient estimation

$$\nabla_{\Theta, \phi} \mathcal{L}(\Theta, \phi) \approx \frac{n}{m} \sum_{i=1}^m \nabla_{\Theta, \phi} \mathcal{L}_i(\Theta, \phi), \quad i \sim \mathbf{U}(1, \dots, n) \tag{39}$$

We need to take gradient from intergral in ELBO, we cant just estimate this integral by Monte-Carlo. Lets make reparameterization trick

$$\begin{aligned}
\int_{z_i} Q(z_i|x_i, \phi) \log p(x_i|z_i, \Theta) dz_i &= \int_{\epsilon} p(\epsilon) \log p(x_i|\epsilon\sigma(x_i, \phi) + \mu(x_i, \phi), \Theta) d\epsilon \approx \\
&\approx \frac{1}{V} \sum_{v=1}^V \log p(x_i|\epsilon_v\sigma(x_i, \phi) + \mu(x_i, \phi), \Theta), \quad \epsilon_v \sim \mathbf{N}(\epsilon_v|0, I) \\
p(x_i|\epsilon_v\sigma(x_i, \phi) + \mu(x_i, \phi), \Theta) &= \mathbf{N}(x_i|f_{\Theta}(\epsilon_v\sigma(x_i, \phi) + \mu(x_i, \phi)), s^2 I)
\end{aligned} \tag{40}$$

Now everything is done, we are ready to write final optimization problem

$$\begin{aligned}
&\nabla_{\Theta, \phi} \mathcal{L}(\Theta, \phi) \\
&\approx \frac{n}{m} \sum_{i=1}^m \nabla_{\Theta, \phi} \left[-\frac{1}{V} \sum_{v=1}^V \|f_{\Theta}(\epsilon_v\sigma(x_i, \phi) + \mu(x_i, \phi)) - x_i\|^2 \right. \\
&\quad \left. + \frac{1}{2} \sum_{j=1}^d (1 + \log \sigma_j(x_i, \phi)^2 - \sigma_j(x_i, \phi)^2 - \mu_j(x_i, \phi)^2) \right]
\end{aligned}$$

V – number of sampling for Monte-Carlo estimation; m – batch size; n – training size;

Θ, ϕ – parameters of NN;

$f_{\Theta}, x_i \in \mathbb{R}^D$ (D – data dimension); $\sigma_{\phi}, \mu_{\phi} \in \mathbb{R}^d$ (d – latent dimension)

(41)

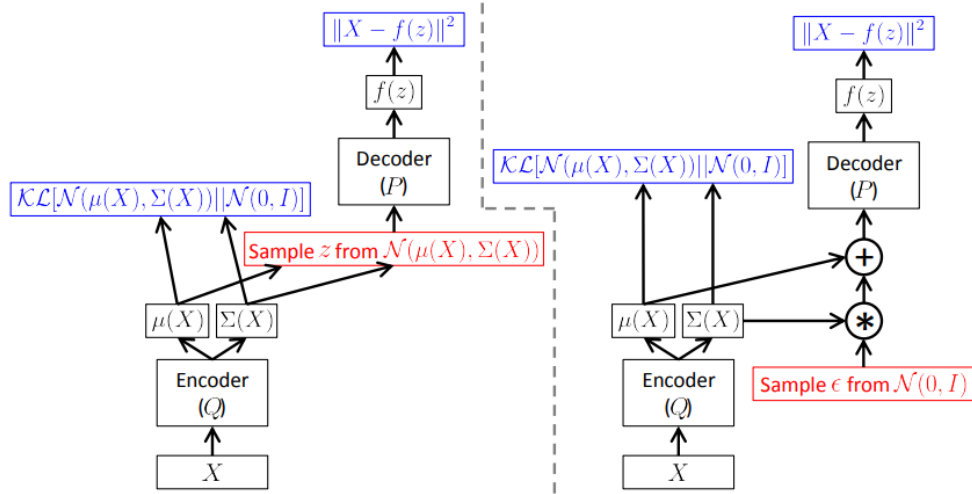


Figure 1: Variational Auto Encoder

This procedure calls double stochastic derivation or variational Bayesian derivation. Schematically it's looks like Figure 1

2.5 Prior hole problem

It is necessary to discuss the main problem of the VAE. As we remember the ELBO looks like

$$\begin{aligned}\mathcal{L}(\Theta, Q) &= \mathbb{E}_{z \sim Q} \log P(X|Z, \Theta) - \mathbb{KL}(Q(Z|X, \Phi) || P(Z)) \\ &= \text{Reconstruction} - \text{Regularization}\end{aligned}\tag{42}$$

Here we have two terms: reconstruction (standard loss between generated and true images like MSE) and regularization (KL divergence between approximation of the posterior and prior). If we drop regularization, we will obtain standard autoencoder, that is, it will not be a generative model. Regularization term is very important because it makes our model generative. But there is some problem called prior hole problem, it is shown in Figure 2

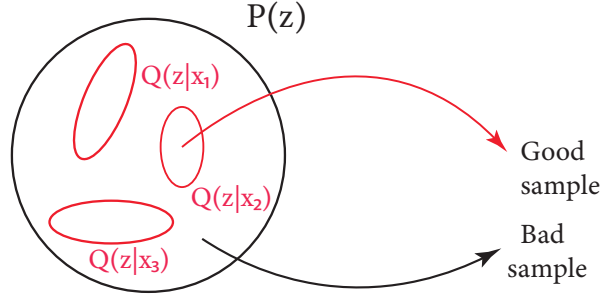


Figure 2: Prior hole problem in VAE. We can take sample from prior distribution that decoder have never seen, so we will obtain bad quality image. In other words, in our prior distribution formed holes.

3 Connection between VAE and EM-algorithm

To provide the connection between VAE and EM-algorithm we have to rewrite ELBO as the following:

$$\begin{aligned}\mathcal{L}_X(\Theta, Q) &= \mathbb{E}_{z \sim Q} \log P(X|Z, \Theta) - \mathbb{KL}(Q(Z|X, \Phi) || P(Z)) = \\ &= \int Q(Z|X, \Phi) \log P(X|Z, \Theta) dZ - \int Q(Z|X, \Phi) \log \frac{Q(Z|X, \Phi)}{P(Z)} dz \\ &= \int Q(Z|X, \Phi) \log \frac{P(X, Z|\Theta)}{Q(Z|X, \Phi)} dZ \\ &= \begin{cases} \log P(X|\Theta) - \mathbb{KL}(Q(Z|X, \Phi) || P(Z|X, \Theta)) \\ \mathbb{E}_Q \log P(X, Z|\Theta) - \mathbb{E}_Q \log Q(Z|X, \Phi) \end{cases}\end{aligned}\tag{43}$$

As we remember, in the EM-algorithm we have to calculate posterior distribution, $P(Z|X, \Theta)$. Unfortunately, in VAE we cannot calculate it, since $P(X|Z, \Theta)$ is a neural network (thus we cannot provide bayesian inference). So, the main difference between VAE and EM-algorithm is the following: **instead of bayesian inference, we approximate the true posterior**. This approximation called *variational inference*.

E-step:

$$\boxed{\begin{aligned} q(z) &= \operatorname{argmax}_Q \mathcal{L}_X(\Theta, Q) = \operatorname{argmin}_\Phi \mathbb{KL}(Q(Z|X, \Phi) || P(Z|X, \Theta)) \\ &= \operatorname{argmin}_\Phi \mathbb{KL}(Q(Z|X, \Phi) || P(Z)) \end{aligned}} \quad (44)$$

M-step:

$$\boxed{\theta_{new} = \operatorname{argmax}_\theta \mathbb{E}_Q \log P(X, Z|\Theta)} \quad (45)$$

4 Information theory

Here some details about information theory needed in generative models will be discussed. First of all, let's consider an entropy of random variable X with distribution $P(X)$ (everything will be considered for discrete variables):

$$\mathbb{H}(X) = - \sum_i P(x_i) \log P(x_i) \quad (46)$$

Here x_i is a realization of the random variable X . This is called a measure of uncertainty. Thus, entropy has a maximum in cases when we do not know anything about variable, so $P = 0.5$. Let's show it in an example of random variable with two possible realizations:

$$\mathbb{H}(X) = -P(x_1) \log P(x_1) - (1 - P(x_1)) \log(1 - P(x_1)) \quad (47)$$

Let's take a derivative (we can't do it because of the discreteness, but who cares:))

$$\frac{\partial \mathbb{H}}{\partial P}(X) = -[\log P(x_1) - \log(1 - P(x_1))] = 0 \quad (48)$$

We have following extremes: $P = 0.5, 0, 1$ and 0.5 is a maximum.

Also, we can calculate a KL-divergence between two probability distributions to measure the dissimilarity between them

$$\mathbb{KL}(p||q) = \sum_i p_i \log \frac{p_i}{q_i} = \sum_i p_i \log p_i - \sum_i p_i \log q_i = -\mathbb{H}(p) + \mathbb{H}(p, q) \quad (49)$$

where $\mathbb{H}(p, q)$ - cross entropy. This can be interpreted as the average number of bits needed to encode sample coming from p when we use samples from q . So, if samples from both distributions are same, then the number of bits is equal. Then the regular entropy is the expected number of bits to encode sample if we use the true model. Thus, the KL-divergence is the average number of extra bits needed to encode the data, due to the fact that we used distribution q to encode the data instead of the true distribution p . If these distributions are same, then extra bits is equal to zero.

References

- [Doe16] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- [KW13] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.