

# Variational AutoEncoder

Shuguo.J

2020 年 10 月 8 日

## Content

AutoEncoder 注重 encoder 模块, 将高维数据转换为 code/representation; 而 Variational AutoEncoder 则注重 decoder 模块, 根据服从某一分布的样本 code 来生成样本。在 AutoEncoder 中, 隐变量  $z$  由 encoder 编码而来, 因此我们并不知道其它服从的分布的形式。这导致我们无法隐变量  $z$  服从的分布中进行采样, 然后通过 decoder 来生成数据样本。Variational AutoEncoder 则在训练的过程中通过 KL Divergence 来约束隐变量  $z$  的分布, 用先验概率  $P(z)$  来拟合复杂的后验概率  $P(z|x)$ 。

在 encoder 模块, Variational AutoEncoder 不是输出固定的 code/representation, 而是输出  $\mu$  和  $\sigma^2$ , 之后在从  $N(\mu, \sigma^2)$  (假设后验概率  $P(z|x)$  为高斯分布) 采样得到 code- $z$ 。为了使得模型可导, 这里引入随机变量  $\epsilon \sim N(0, 1)$ , 令  $z = \mu + \epsilon \times \sigma$ , 并以此来代替上述的随机采样。在这里, 损失函数就能够对  $\mu$  和  $\sigma^2$  进行求导, 而其无需对  $\epsilon$  求导。同时,  $z \sim N(\mu, \sigma^2)$ 。

**Variational AutoEncoder** 的完整公式推导我还没有完全理解好。  
待更新...

## Appendix

KL divergence 如下所示:

$$D_{kl}(p(x)|q(x)) = \int p(x) \log \frac{p(x)}{q(x)} dx \quad (1)$$

KL 散度代表着用分布  $q(x)$  来拟合  $p(x)$  的差异程度。

当  $p(x)$  和  $q(x)$  均为高斯分数时,

$$\begin{aligned}
D_{kl}(p(x)|q(x)) &= \int p(x) \times \left( \log \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} - \log \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} \right) dx \\
&= \int p(x) \times \left( \log \frac{1}{\sqrt{2\pi\sigma_1^2}} + \log e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} - \log \frac{1}{\sqrt{2\pi\sigma_2^2}} - \log e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} \right) dx \\
&= \int p(x) \times \left( -\frac{1}{2} \log 2\pi - \log \sigma_1 - \frac{(x-\mu_1)^2}{2\sigma_1^2} + \frac{1}{2} \log 2\pi + \log \sigma_2 + \frac{(x-\mu_2)^2}{2\sigma_2^2} \right) dx \\
&= \int p(x) \times \left( \log \frac{\sigma_2}{\sigma_1} - \frac{(x-\mu_1)^2}{2\sigma_1^2} + \frac{(x-\mu_2)^2}{2\sigma_2^2} \right) dx \\
&= \log \frac{\sigma_2}{\sigma_1} + \int p(x) \times \frac{(x-\mu_2)^2}{2\sigma_2^2} dx - \int p(x) \frac{(x-\mu_1)^2}{2\sigma_1^2} dx \\
&= \log \frac{\sigma_2}{\sigma_1} + \frac{1}{2\sigma_2^2} \int p(x) (x - \mu_1 + \mu_1 - \mu_2)^2 dx - \frac{1}{2} \\
&= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma^2 + (\mu_2 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}
\end{aligned}$$

(2)