

# KL divergence between two multivariate Gaussians

jojonki

May 2018

## 1 導出するよ

VAEでは多次元正規分布 $\mathcal{N}(\mathbf{0}, \mathbf{I})$ と $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$ のKLダイバージェンスを求めたが、 $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\sigma}_1^2)$ と $\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\sigma}_2^2)$ のKLダイバージェンスを求めたいと思う。それぞれ $q(z)$ ,  $p(z)$ として考えてみる。

$$\begin{aligned} D_{KL}(q(z)||p(z)) &= \int q(z) \log \frac{q(z)}{p(z)} dz \\ &= \int q(z) \left\{ \log q(z) - \log p(z) \right\} dz \\ &= \int q(z) \log q(z) dz - \int q(z) \log p(z) dz \end{aligned} \tag{1}$$

第1項と第2項でそれぞれ計算する。第1項は以前にVAEの計算のときに求めた方法と同じ。第2項が異なる分布間（4変数）なので少しだけトリッキー。

第1項

$$\begin{aligned}
\int q(z) \log q(z) dz &= \int \mathcal{N}(z; \mu_1, \sigma_1^2) \log \mathcal{N}(z; \mu_1, \sigma_1^2) dz \\
&= \sum_j^J E_{q(z_j)} \left[ \log \frac{1}{\sqrt{2\pi\sigma_{1,j}^2}} \exp \left( -\frac{(z_j - \mu_{1,j})^2}{2\sigma_{1,j}^2} \right) \right] \\
&= \sum_j^J E_{q(z_j)} \left[ -\frac{1}{2} \log(2\pi\sigma_{1,j}^2) - \frac{(z_j - \mu_{1,j})^2}{2\sigma_{1,j}^2} \right] \\
&= \sum_j^J \left\{ -\frac{1}{2} \log(2\pi\sigma_{1,j}^2) - E_{q(z_j)} \left[ \frac{(z_j - \mu_{1,j})^2}{2\sigma_{1,j}^2} \right] \right\} \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{1,j}^2 - \frac{1}{2} \sum_{j=1}^J E_{q(z_j)} \left[ \frac{(z_j - \mu_{1,j})^2}{\sigma_{1,j}^2} \right] \tag{2} \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{1,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{1,j}^2} E_{q(z_j)} \left[ (z_j - \mu_{1,j})^2 \right]
\end{aligned}$$

第3項の右側は分散の定義そのもの.

$$\begin{aligned}
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{1,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{1,j}^2} \sigma_{1,j}^2 \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\log \sigma_{1,j}^2 + 1)
\end{aligned}$$

第2項

$$\begin{aligned}
\int q(z) \log p(z) dz &= \int \mathcal{N}(z; \mu_1, \sigma_1^2) \log \mathcal{N}(z; \mu_2, \sigma_2^2) dz \\
&= \sum_j^J E_{q(z_j)} \left[ \log \frac{1}{\sqrt{2\pi\sigma_{2,j}^2}} \exp \left( -\frac{(z_j - \mu_{2,j})^2}{2\sigma_{2,j}^2} \right) \right] \\
&= \sum_j^J E_{q(z_j)} \left[ -\frac{1}{2} \log(2\pi\sigma_{2,j}^2) - \frac{(z_j - \mu_{2,j})^2}{2\sigma_{2,j}^2} \right] \\
&= \sum_j^J \left\{ -\frac{1}{2} \log(2\pi\sigma_{2,j}^2) - E_{q(z_j)} \left[ \frac{(z_j - \mu_{2,j})^2}{2\sigma_{2,j}^2} \right] \right\} \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J E_{q(z_j)} \left[ \frac{(z_j - \mu_{2,j})^2}{\sigma_{2,j}^2} \right] \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{2,j}^2} E_{q(z_j)} \left[ (z_j - \mu_{2,j})^2 \right] \\
&\quad \text{ここから先程と違う．第3項を展開していく．} \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{2,j}^2} E_{q(z_j)} \left[ z_j^2 - 2z_j\mu_{2,j} + \mu_{2,j}^2 \right] \tag{3} \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{2,j}^2} \left[ E_{q(z_j)} z_j^2 - 2E_{q(z_j)} z_j \mu_{2,j} + E_{q(z_j)} \mu_{2,j}^2 \right] \\
&\quad q(z) \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\sigma}_1^2) \text{に関する期待値を}\langle X \rangle \text{の記法で便宜上書く．} \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{2,j}^2} \left[ \langle z_j^2 \rangle - 2\langle z_j \rangle \mu_{2,j} + \mu_{2,j}^2 \right] \\
&\quad \text{分散の公式から}\langle z_j^2 \rangle = \sigma_{1,j}^2 + \mu_{1,j}^2 \text{であり，}\langle z_j \rangle = \mu_j \text{であるから，} \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{2,j}^2} \left[ \sigma_{1,j}^2 + \mu_{1,j}^2 - 2\mu_{1,j}\mu_{2,j} + \mu_{2,j}^2 \right] \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{1}{\sigma_{2,j}^2} \left[ \sigma_{1,j}^2 + (\mu_{1,j} - \mu_{2,j})^2 \right] \\
&= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \left[ \frac{\sigma_{1,j}^2}{\sigma_{2,j}^2} + \frac{(\mu_{1,j} - \mu_{2,j})^2}{\sigma_{2,j}^2} \right]
\end{aligned}$$

まとめ

第1項から第2項を引けば良いので,

$$\begin{aligned}
D_{KL}(q(z)||p(z)) &= \int q(z) \log \frac{q(z)}{p(z)} dz \\
&= \int q(z) \{ \log q(z) - \log p(z) \} dz \\
&= \int q(z) \log q(z) dz - \int q(z) \log p(z) dz \\
&= \left\{ -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\log \sigma_{1,j}^2 + 1) \right\} - \left\{ -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \sigma_{2,j}^2 - \frac{1}{2} \sum_{j=1}^J \left[ \frac{\sigma_{1,j}^2}{\sigma_{2,j}^2} + \frac{(\mu_{1,j} - \mu_{2,j})^2}{\sigma_{2,j}^2} \right] \right\} \\
&= -\frac{1}{2} \sum_{j=1}^J \left[ \log \frac{\sigma_{1,j}^2}{\sigma_{2,j}^2} - \frac{\sigma_{1,j}^2}{\sigma_{2,j}^2} - \frac{(\mu_{1,j} - \mu_{2,j})^2}{\sigma_{2,j}^2} + 1 \right]
\end{aligned} \tag{4}$$

ちなみにVAEのときを復習すると,  $p(z)$ はVAEの設定から  $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  であり,  $q(z|x)$ は  $\sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$  となる. これは今回求めたものに,  $\mu_2 = 0, \sigma_2 = 1$ を代入すれば同じ値になることが確認できる.

$$\begin{aligned}
D_{KL}(q(z|x)||p(z)) &= \int q(z|x) \log \frac{q(z|x)}{p(z)} dz \\
&= \int q(z|x) \{ \log q(z|x) - \log p(z) \} dz \\
&= \int q(z|x) \log q(z|x) dz - \int q(z|x) \log p(z) dz \\
&= \left( -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\log \sigma_j^2 + 1) \right) - \left( -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\sigma_j^2 + \mu_j^2) \right) \\
&= -\frac{1}{2} \sum_{j=1}^J \left( 1 + \log \sigma^2 - \sigma_j^2 - \mu_j^2 \right)
\end{aligned} \tag{5}$$

## 2 参考

- <https://stats.stackexchange.com/questions/7440/kl-divergence-between-two-univariate-gaussians>