# PAC - Bayes

## Question 1

Consider two multivariate Gaussian distributions over $\mathbb{R}^r$ - $N(\mu_0, \Sigma_0)$ and $N(\mu_1, \Sigma_1)$ where $\Sigma_0$ and $\Sigma_1$ are non-singular (positive definite). Then it holds that:

$$KL(N(\mu_0, \Sigma_0) \| N(\mu_1, \Sigma_1)) = \frac{1}{2}\left[ tr(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1}(\mu_1 - \mu_0) - r + \ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) \right]$$

## Solution

First, let's sign the two distributions as:

$$P = N(\mu_0, \Sigma_0) \quad \text{and} \quad Q = N(\mu_1, \Sigma_1).$$

By definition: $KL(P \| Q) = \mathbb{E}_P\left[ \ln\left(\frac{P}{Q}\right) \right].$

Remember that the density function of a multivariate Gaussian distribution $N(\mu, \Sigma)$ is:

$$p(x) = \frac{1}{(2\pi)^{\frac{r}{2}} (\det \Sigma)^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$\Rightarrow$

$$\frac{P(x)}{Q(x)} = P(x) \cdot \frac{1}{Q(x)} = \frac{1}{(2\pi)^{\frac{r}{2}}(\det\Sigma_0)^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(x-\mu_0)^T \Sigma_0^{-1}(x-\mu_0)} \cdot \frac{(2\pi)^{\frac{r}{2}} (\det\Sigma_1)^{\frac{1}{2}}}{e^{-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1)}}$$

$$= \frac{(\det\Sigma_1)^{\frac{1}{2}}}{(\det\Sigma_0)^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(x-\mu_0)^T \Sigma_0^{-1}(x-\mu_0) - (-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1))}$$

$$= \left(\frac{\det\Sigma_1}{\det\Sigma_0}\right)^{\frac{1}{2}} \cdot e^{\frac{1}{2}\left[(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1) - (x-\mu_0)^T \Sigma_0^{-1}(x-\mu_0)\right]}$$

Thus:

$$KL(P \| Q) = \mathbb{E}_p\left[ \ln\left(\frac{P}{Q}\right) \right]$$

$$= \mathbb{E}_p\left[ \ln\left( \left(\frac{\det \Sigma_1}{\det \Sigma_0}\right)^{\frac{1}{2}} \cdot e^{\frac{1}{2}\left[ (x-M_1)^T \Sigma_1^{-1}(x-M_1) - (x-M_0)^T \Sigma_0^{-1}(x-M_0) \right]} \right) \right]$$

$$\left[ \ln(a \cdot b) = \ln(a) + \ln(b) \right] = \mathbb{E}_p\left[ \ln\left[ \left(\frac{\det \Sigma_1}{\det \Sigma_0}\right)^{\frac{1}{2}} \right] + \ln\left[ e^{\frac{1}{2}\left[ (x-M_1)^T \Sigma_1^{-1}(x-M_1) - (x-M_0)^T \Sigma_0^{-1}(x-M_0) \right]} \right] \right]$$

$$= \mathbb{E}_p\left[ \frac{1}{2} \ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \frac{1}{2}\cdot\left( (x-M_1)^T \Sigma_1^{-1}(x-M_1) - (x-M_0)^T \Sigma_0^{-1}(x-M_0) \right) \right]$$

linearity of $\mathbb{E}(\cdot)$ $\searrow$

$$= \mathbb{E}_p\left[ \frac{1}{2}\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) \right] + \mathbb{E}_p\left[ \frac{1}{2}(x-M_1)^T \Sigma_1^{-1}(x-M_1) \right]$$

$$- \mathbb{E}_p\left[ \frac{1}{2}(x-M_0)^T \Sigma_0^{-1}(x-M_0) \right]$$

$$= \frac{1}{2}\cdot\left[ \ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \mathbb{E}_p\left[ (x-M_1)^T \Sigma_1^{-1}(x-M_1) \right] - \mathbb{E}_p\left[ (x-M_0)^T \Sigma_0^{-1}(x-M_0) \right] \right]$$

### claim

For each $x \in \mathbb{R}^k$, $A \in \mathbb{R}^{k \times r}$: $\quad x^T A x = tr(A x x^T)$

### proof

$$x^T A x = (x_1 \cdots x_r)\begin{pmatrix} A^1 x \\ \vdots \\ A^k x \end{pmatrix} = (x_1 \cdots x_r)\begin{pmatrix} (Ax)_1 \\ \vdots \\ (Ax)_k \end{pmatrix} = (x_1 \cdots x_r)\begin{pmatrix} \sum_{j=1}^{r} a_{1j}x_j \\ \vdots \\ \sum_{j=1}^{r} a_{rj}x_j \end{pmatrix} =$$

$$= \sum_{j=1}^{k} a_{1j}x_1 x_j + \cdots + \sum_{j=1}^{k} a_{kj}x_k x_j =$$

$$= \sum_{j=1}^{r} a_{1j}(xx^T)_{j1} + \cdots + \sum_{j=1}^{h} a_{rj}(xx^T)_{jr} =$$

$$= (Axx^T)_{11} + \cdots + (Axx^T)_{rr} = tr(Axx^T)$$

$\square$

Thus:

$$KL(P||Q) = \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \mathbb{E}_P[(x-M_1)^T \Sigma_1^{-1}(x-M_1)] - \mathbb{E}_P[(x-M_0)^T \Sigma_0^{-1}(x-M_0)]\right]$$

$$= \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \mathbb{E}_P[\text{tr}(\Sigma_1^{-1}(x-M_1)(x-M_1)^T)] - \mathbb{E}_P[\text{tr}(\Sigma_0^{-1}(x-M_0)(x-M_0)^T)]\right]$$

tr(E[A]) = E[tr(A)]
by linearity
of E[·]

$$= \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \text{tr}(\mathbb{E}_P[\Sigma_1^{-1}(x-M_1)(x-M_1)^T]) - \text{tr}(\mathbb{E}_P[\Sigma_0^{-1}(x-M_0)(x-M_0)^T])\right]$$

$$= \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \text{tr}(\Sigma_1^{-1}\mathbb{E}_P[xx^T - M_1 x^T - x M_1^T + M_1 M_1^T])\right.$$
$$\left. - \text{tr}(\Sigma_0^{-1}\mathbb{E}_P[(x-M_0)(x-M_0)^T])\right] \quad (=)$$

By definition, $P = N(M_0, \Sigma_0)$ thus

$$\mathbb{E}_P[x] = M_0 \quad \text{and} \quad \mathbb{E}_P[(x-M_0)(x-M_0)^T] = \Sigma_0$$

$$\text{and} \quad \mathbb{E}_P[xx^T] = \Sigma_0 + M_0 M_0^T, \text{ so:}$$

$$(=) \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \text{tr}(\Sigma_1^{-1}(\Sigma_0 + M_0 M_0^T - M_1 M_0^T - M_0 M_1^T + M_1 M_1^T))\right.$$
$$\left. - \text{tr}(\Sigma_0^{-1}\Sigma_0)\right]$$

$$= \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) + \text{tr}(\Sigma_1^{-1}\Sigma_0) + \text{tr}(\Sigma_1^{-1}(M_0 M_0^T - M_1 M_0^T - M_0 M_1^T + M_1 M_1^T))\right.$$
$$\left. - \text{tr}(I_r)\right]$$

$X^T A X$
$= A x x^T$
by claim
we proved

$$= \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) - r + \text{tr}(\Sigma_1^{-1}\Sigma_0) + \text{tr}(M_0^T \Sigma_1^{-1} M_0 - M_0^T \Sigma_1^{-1} M_1 - M_1^T \Sigma_1^{-1} M_0 + M_1^T \Sigma_1^{-1} M_1)\right]$$

$$= \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) - r + \text{tr}(\Sigma_1^{-1}\Sigma_0) + (M_1 - M_0)^T \Sigma_1^{-1}(M_1 - M_0)\right]$$

$$\Rightarrow$$

$$KL(N(M_0, \Sigma_0) || N(M_1, \Sigma_1)) = \frac{1}{2}\left[\ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) - r + \text{tr}(\Sigma_1^{-1}\Sigma_0) + (M_1 - M_0)^T \Sigma_1^{-1}(M_1 - M_0)\right]$$

as required. □