

1 KL divergence

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$

2 Expected value

$$\mathbb{E}[X] = \int x f(x) dx$$

- **Linearity**

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

$$\mathbb{E}[aX] = a \mathbb{E}[X]$$

3 Multivariate Normal Distribution

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) \quad (1)$$

- **log likelihood**

$$\begin{aligned} \ln(p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})) &= -\frac{k \ln(2\pi)}{2} - \frac{\ln(|\boldsymbol{\Sigma}|)}{2} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= -\frac{1}{2} \left[\ln \left((2\pi)^k |\boldsymbol{\Sigma}| \right) + (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \end{aligned}$$

- **KL divergence**

$$\begin{aligned} D_{KL}(\mathcal{N}_a || \mathcal{N}_b) &= \frac{1}{2} \left(\text{Tr}(\boldsymbol{\Sigma}_b^{-1} \boldsymbol{\Sigma}_a) + (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a)^T \boldsymbol{\Sigma}_b^{-1} (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a) - k + \log \frac{|\boldsymbol{\Sigma}_b|}{|\boldsymbol{\Sigma}_a|} \right) \\ &= \frac{1}{2} \left(\|\mathbf{L}_b \setminus \mathbf{L}_a\|_F^2 + \|\mathbf{L}_b \setminus (\boldsymbol{\mu}_b - \boldsymbol{\mu}_a)\|_2^2 - k + 2 \left(\log \text{diag}(\mathbf{L}_b)^T \mathbf{1} - \log \text{diag}(\mathbf{L}_a)^T \mathbf{1} \right) \right) \end{aligned}$$

- **Affine Transform**

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$$

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu}_x + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}_x \mathbf{A}^T)$$

- **Linear Gaussian systems** Given a linear system:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$$

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \boldsymbol{\Sigma}_y)$$

We have the following:

$$\begin{aligned}
p(\mathbf{x} \mid \mathbf{y}) &= \mathcal{N} \left(\mathbf{x} \mid \boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y} \right) \\
\boldsymbol{\mu}_{x|y} &= \boldsymbol{\Sigma}_{x|y} \left[\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \mathbf{b}) + \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\mu}_x \right] \\
\boldsymbol{\Sigma}_{x|y} &= \boldsymbol{\Sigma}_x^{-1} + \mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A} \\
p(\mathbf{y}) &= \mathcal{N} \left(\mathbf{y} \mid \mathbf{A} \boldsymbol{\mu}_x + \mathbf{b}, \boldsymbol{\Sigma}_y + \mathbf{A} \boldsymbol{\Sigma}_x \mathbf{A}^T \right)
\end{aligned}$$

- **quadratic relations**

$$- \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\mathbf{x}^T \mathbf{A} \mathbf{x}] = \text{Tr}(\mathbf{A} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$$

-

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{b}, \mathbf{B})} \left[(\mathbf{a} - \mathbf{A} \mathbf{x})^T \boldsymbol{\Sigma}^{-1} (\mathbf{a} - \mathbf{A} \mathbf{x}) \right] = (\mathbf{a} - \mathbf{A} \mathbf{b})^T \boldsymbol{\Sigma}^{-1} (\mathbf{a} - \mathbf{A} \mathbf{b}) + \text{Tr} \left(\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A} \mathbf{B} \right)$$

4 Gamma distribution

$x \sim Ga(a, b)$ where a is called the shape and b the rate.

$$Ga(x|a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} \exp(-bx)$$

- $\mathbb{E}[x] = \frac{a}{b}$
- $\mathbb{E}_{x \sim Ga(a, b)} [\ln x] = \psi(a) - \ln(b)$
where ψ is the polygamma function.

4.1 Inverse gamma

If $x \sim Ga(a, b)$ and $y = \frac{1}{x}$, then $y \sim IG(a, b)$

$$IG(y|a, b) = \frac{b^a}{\Gamma(a)} y^{-(a+1)} \exp(-b/y)$$

4.2 Inverse Wishart (IW)

This distribution is used in Bayesian statistics as the conjugate prior for the covariance matrix of a multivariate normal distribution:

$$\begin{aligned}
\boldsymbol{\Sigma} &\sim IW(\mathbf{S}^{-1}, v + D + 1) \\
IW(\boldsymbol{\Sigma} \mid \mathbf{S}, v) &= \frac{1}{Z(\mathbf{S}, v)} |\boldsymbol{\Sigma}|^{(v+D+1)/2} \exp \left(-\frac{1}{2} \text{Tr}(\mathbf{S}^{-1} \boldsymbol{\Sigma}^{-1}) \right) \\
Z(\mathbf{S}, v) &= |\mathbf{S}|^{-v/2} 2^{vD/2} \Gamma_D v / 2
\end{aligned}$$

where $\mathbf{S} \succ 0$