

### 2.3.5 多次元ガウス分布

分布関数

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (2.72)$$

( $\boldsymbol{\Sigma}$  は正定値対称行列)

規格化定数の求め方

$$C = \int d\mathbf{x} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (\text{A.5.1})$$

を考える。 $\boldsymbol{\Sigma}$  が正定値であるため、 $\boldsymbol{\Sigma}^{-1}$  も正定値だから、下三角行列  $\mathbf{L}$  を使ってコレスキー分解

$$\boldsymbol{\Sigma}^{-1} = \mathbf{L} \mathbf{L}^T \quad (\text{A.5.2})$$

できる。積分変数の変数変換

$$\mathbf{y} = \mathbf{L}^T (\mathbf{x} - \boldsymbol{\mu}) \quad (\text{A.5.3})$$

とすると、

$$d\mathbf{x} = \frac{1}{|\mathbf{L}^T|} d\mathbf{y} = \sqrt{|\boldsymbol{\Sigma}|} d\mathbf{y} \quad (\text{A.5.4})$$

より

$$\begin{aligned} C &= \sqrt{|\boldsymbol{\Sigma}|} \int d\mathbf{y} \exp \left\{ -\frac{1}{2} \mathbf{y}^T \mathbf{y} \right\} \\ &= \sqrt{|\boldsymbol{\Sigma}|} \sqrt{(2\pi)^D} \\ &= \sqrt{(2\pi)^D |\boldsymbol{\Sigma}|} \end{aligned} \quad (\text{A.5.5})$$

対数表示

$$\ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{1}{2} \left\{ (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) + \ln |\boldsymbol{\Sigma}| + D \ln 2\pi \right\} \quad (2.73)$$

$\Sigma$  が対角行列の場合

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_D^2 \end{pmatrix} \quad (2.74)$$

を代入すると

$$\begin{aligned} \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) &= -\frac{1}{2} \sum_{d=1}^D \left\{ \frac{(x_d - \mu_d)^2}{\sigma_d^2} + \ln \sigma_d^2 + \ln 2\pi \right\} \\ &= \ln \prod_{d=1}^D \mathcal{N}(x_d|\mu_d, \sigma_d^2) \end{aligned} \quad (2.75)$$

となる。ここで、右辺は1次元ガウス分布の式

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\} \quad (2.64)$$

で、計算には以下を使った。

$$\begin{aligned} \Sigma^{-1} &= \begin{pmatrix} \frac{1}{\sigma_1^2} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\sigma_D^2} \end{pmatrix} \\ (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) &= \sum_{d=1}^D \frac{(x_d - \mu_d)^2}{\sigma_d^2} \\ \ln |\Sigma| &= \sum_{d=1}^D \ln \sigma_d^2 \end{aligned}$$

基本的な期待値

$$\begin{aligned} \langle \mathbf{x} \rangle &= \int d\mathbf{x} \mathbf{x} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) \\ &= \int d\mathbf{x} \mathbf{x} \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= \int d\tilde{\mathbf{x}} (\tilde{\mathbf{x}} + \boldsymbol{\mu}) \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} \tilde{\mathbf{x}}^T \Sigma^{-1} \tilde{\mathbf{x}} \right\} \\ &= \boldsymbol{\mu} \int d\tilde{\mathbf{x}} \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} \tilde{\mathbf{x}}^T \Sigma^{-1} \tilde{\mathbf{x}} \right\} \\ &= \boldsymbol{\mu} \end{aligned} \quad (2.76)$$

$$\begin{aligned}
\langle \mathbf{x} \mathbf{x}^T \rangle &= \int d\mathbf{x} \mathbf{x} \mathbf{x}^T \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\
&= \int d\mathbf{x} \mathbf{x} \mathbf{x}^T \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\
&= \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \int d\tilde{\mathbf{x}} (\tilde{\mathbf{x}} + \boldsymbol{\mu}) (\tilde{\mathbf{x}} + \boldsymbol{\mu})^T \exp \left\{ -\frac{1}{2} \tilde{\mathbf{x}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}} \right\} \\
&= \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \int d\tilde{\mathbf{x}} (\tilde{\mathbf{x}} \tilde{\mathbf{x}}^T + \boldsymbol{\mu} \tilde{\mathbf{x}}^T + \tilde{\mathbf{x}} \boldsymbol{\mu}^T + \boldsymbol{\mu} \boldsymbol{\mu}^T) \\
&\quad \exp \left\{ -\frac{1}{2} \tilde{\mathbf{x}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}} \right\} \\
&= \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \int d\tilde{\mathbf{x}} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T \exp \left\{ -\frac{1}{2} \tilde{\mathbf{x}}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}} \right\} + \boldsymbol{\mu} \boldsymbol{\mu}^T
\end{aligned}$$

ここで、右辺第一項の  $i, j$  成分を  $I_{ij}$  として、改めて  $\tilde{\mathbf{x}}$  を  $\mathbf{x}$  と書くと

$$I_{ij} = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \int d\mathbf{x} x_i x_j \exp \left\{ -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right\}$$

ここで、 $\boldsymbol{\Sigma}$  の逆行列を  $\mathbf{M} = \boldsymbol{\Sigma}^{-1}$  とすると、

$$\begin{aligned}
I_{ij} &= \sqrt{\frac{|\mathbf{M}|}{(2\pi)^D}} \int d\mathbf{x} x_i x_j \exp \left\{ -\frac{1}{2} \mathbf{x}^T \mathbf{M} \mathbf{x} \right\} \\
&= \sqrt{\frac{|\mathbf{M}|}{(2\pi)^D}} (-2) \frac{\partial}{\partial M_{ij}} \int d\mathbf{x} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \mathbf{M} \mathbf{x} \right\} \\
&= \sqrt{\frac{|\mathbf{M}|}{(2\pi)^D}} (-2) \frac{\partial}{\partial M_{ij}} \sqrt{\frac{(2\pi)^D}{|\mathbf{M}|}} \\
&= -2 |\mathbf{M}|^{\frac{1}{2}} \frac{\partial}{\partial M_{ij}} |\mathbf{M}|^{-\frac{1}{2}} \\
&= \frac{1}{|\mathbf{M}|} \frac{\partial |\mathbf{M}|}{\partial M_{ij}} \\
&= \frac{\tilde{M}_{ij}^T}{|\mathbf{M}|} \\
&= \Sigma_{ij}
\end{aligned}$$

ここで、 $\tilde{M}$  は  $\mathbf{M}$  の余因子行列。

以上により

$$\langle \mathbf{x} \mathbf{x}^T \rangle = \tilde{\boldsymbol{\mu}} \boldsymbol{\mu}^T + \boldsymbol{\Sigma} \quad (2.77)$$

エントロピー

$$\begin{aligned}
H[\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})] &= -\langle \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \rangle \\
&= \frac{1}{2} \left\{ \langle (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \rangle + \ln |\boldsymbol{\Sigma}| + D \ln 2\pi \right\}
\end{aligned} \tag{2.78}$$

右辺中括弧内第一項は計算できる。(テキストとは2行目が違う?)

$$\begin{aligned}
&\langle (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \rangle \\
&= \langle \text{Tr}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}] \rangle \\
&= \langle \text{Tr}[(\mathbf{x}\mathbf{x}^T - \mathbf{x}\boldsymbol{\mu}^T - \boldsymbol{\mu}\mathbf{x}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T)\boldsymbol{\Sigma}^{-1}] \rangle \\
&= \text{Tr}[(\langle \mathbf{x}\mathbf{x}^T \rangle - \langle \mathbf{x} \rangle \boldsymbol{\mu}^T - \boldsymbol{\mu} \langle \mathbf{x}^T \rangle + \boldsymbol{\mu}\boldsymbol{\mu}^T)\boldsymbol{\Sigma}^{-1}] \\
&= \text{Tr}[(\boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma} - \boldsymbol{\mu}\boldsymbol{\mu}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T)\boldsymbol{\Sigma}^{-1}] \\
&= \text{Tr}[\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}] \\
&= \text{Tr}(\mathbf{I}_D) \\
&= D
\end{aligned} \tag{2.79}$$

(2.78) に (2.79) を代入して

$$H[\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})] = \frac{1}{2} \{ \ln |\boldsymbol{\Sigma}| + D(\ln 2\pi + 1) \} \tag{2.80}$$

KL ダイバージェンス

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \tag{A.5.6}$$

$$q(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) \tag{A.5.7}$$

として、KL ダイバージェンスの定義式 (2.13) を使うと

$$\begin{aligned}
&\text{KL}[q(\mathbf{x})][p(\mathbf{x})] \\
&= -H[\mathcal{N}(\mathbf{x}|\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})] - \langle \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \rangle_{q(\mathbf{x})} \\
&= -\frac{1}{2} \left\{ \ln |\hat{\boldsymbol{\Sigma}}| + D(\ln 2\pi + 1) \right\} \\
&\quad + \frac{1}{2} \left\{ \langle (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \rangle_{q(\mathbf{x})} + \ln |\boldsymbol{\Sigma}| + D \ln 2\pi \right\} \\
&= \frac{1}{2} \left\{ \langle (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \rangle_{q(\mathbf{x})} + \ln \frac{|\boldsymbol{\Sigma}|}{|\hat{\boldsymbol{\Sigma}}|} - D \right\}
\end{aligned} \tag{A.5.8}$$

ここで、中括弧内第一項を考える。

$$\begin{aligned}
& \langle (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \rangle_{q(\mathbf{x})} \\
&= \langle \text{Tr} [(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}] \rangle_{q(\mathbf{x})} \\
&= \langle \text{Tr} [(\mathbf{x}\mathbf{x}^T - \mathbf{x}\boldsymbol{\mu}^T - \boldsymbol{\mu}\mathbf{x}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T) \boldsymbol{\Sigma}^{-1}] \rangle_{q(\mathbf{x})} \\
&= \text{Tr} [\langle \mathbf{x}\mathbf{x}^T \rangle_{q(\mathbf{x})} - \langle \mathbf{x} \rangle_{q(\mathbf{x})} \boldsymbol{\mu}^T - \boldsymbol{\mu} \langle \mathbf{x}^T \rangle_{q(\mathbf{x})} + \boldsymbol{\mu}\boldsymbol{\mu}^T] \boldsymbol{\Sigma}^{-1}] \\
&= \text{Tr} \left[ \left( \hat{\boldsymbol{\mu}}\hat{\boldsymbol{\mu}}^T + \hat{\boldsymbol{\Sigma}} - \hat{\boldsymbol{\mu}}\boldsymbol{\mu}^T - \boldsymbol{\mu}\hat{\boldsymbol{\mu}}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T \right) \boldsymbol{\Sigma}^{-1} \right] \\
&= \text{Tr} \left[ \left( \hat{\boldsymbol{\mu}}\hat{\boldsymbol{\mu}}^T - \hat{\boldsymbol{\mu}}\boldsymbol{\mu}^T - \boldsymbol{\mu}\hat{\boldsymbol{\mu}}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T + \hat{\boldsymbol{\Sigma}} \right) \boldsymbol{\Sigma}^{-1} \right] \\
&= \text{Tr} \left[ \left( (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu})(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu})^T + \hat{\boldsymbol{\Sigma}} \right) \boldsymbol{\Sigma}^{-1} \right] \tag{A.5.9}
\end{aligned}$$

よって、KL ダイバージェンスは

$$\begin{aligned}
& \text{KL}[q(\mathbf{x})][p(\mathbf{x})] \\
&= \frac{1}{2} \left\{ \text{Tr} \left[ \left( (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu})(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu})^T + \hat{\boldsymbol{\Sigma}} \right) \boldsymbol{\Sigma}^{-1} \right] + \ln \frac{|\boldsymbol{\Sigma}|}{|\hat{\boldsymbol{\Sigma}}|} - D \right\} \tag{2.84}
\end{aligned}$$

となる。