

### 演習問題 2.13

2 つのガウス分布  $p(\mathbf{x}) = N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  と  $q(\mathbf{x}) = N(\mathbf{x} | \mathbf{m}, \mathbf{L})$  の間のカルバック – ライブラーダイバージェンス ( 1.13 ) を求めよ。

#### [ カルバック – ライブラーダイバージェンス ]

分布  $p(\mathbf{x})$  と  $q(\mathbf{x})$  の間の相対エントロピーを表す。

$$\begin{aligned} \text{KL}(p \parallel q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left( - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \end{aligned} \quad \cdots ( 1.113 )$$

カルバック – ライブラーダイバージェンスを 2 つの分布  $q(\mathbf{x})$  と  $p(\mathbf{x})$  の間の隔たりを表す尺度として解釈できる。

#### [ 多変量ガウス分布 ]

$$N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad \cdots ( 2.43 )$$

ただし、 $\boldsymbol{\mu}$  は  $D$  次元ベクトル、 $\boldsymbol{\Sigma}$  は  $D \times D$  の共分散行列、 $|\boldsymbol{\Sigma}|$  は  $\boldsymbol{\Sigma}$  の行列式を表す。

#### [ 多変量ガウス分布の正規化 ]

$$\int p(\mathbf{y}) d\mathbf{y} = \prod_{j=1}^D \int_{-\infty}^{\infty} \frac{1}{(2\pi\lambda_j)^{1/2}} \exp \left\{ - \frac{y_j^2}{2\lambda_j} \right\} dy_j = 1 \quad \cdots ( 2.57 )$$

#### [ 多変量ガウス分布の期待値 ( 一次のモーメント ) ]

$$E[\mathbf{x}] = \boldsymbol{\mu} \quad \cdots ( 2.59 )$$

#### [ 多変量ガウス分布の期待値 ( 二次のモーメント ) ]

$$E[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma} \quad \cdots ( 2.62 )$$

[ 解 ]

式 ( 1.113 ) より、2 つのガウス分布  $p(\mathbf{x}) = N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  と  $q(\mathbf{x}) = N(\mathbf{x} | \mathbf{m}, \mathbf{L})$  の間のカルバック – ライブラーダイバージェンス  $\text{KL}(p \parallel q)$  は、

$$\begin{aligned} \text{KL}(p \parallel q) &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \\ &= - \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \ln \left\{ \frac{N(\mathbf{x} | \mathbf{m}, \mathbf{L})}{N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})} \right\} d\mathbf{x} \end{aligned}$$

$$\begin{aligned}
&= - \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \ln \left\{ \frac{\frac{1}{|\mathbf{L}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}) \right\}}{\frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}} \right\} d\mathbf{x} \\
&= - \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \left\{ \frac{1}{2} \ln \frac{|\boldsymbol{\Sigma}|}{|\mathbf{L}|} - \frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}) \right. \\
&\quad \left. + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} d\mathbf{x} \\
&= - \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \ln \frac{|\boldsymbol{\Sigma}|}{|\mathbf{L}|} d\mathbf{x} + \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x} - \mathbf{m})^T \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}) d\mathbf{x} \\
&\quad - \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) d\mathbf{x}
\end{aligned}$$

… ※

と展開できる。ここで、上記の式の第一項については、式 ( 2.57 ) より、

$$\begin{aligned}
&- \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \ln \frac{|\boldsymbol{\Sigma}|}{|\mathbf{L}|} d\mathbf{x} \\
&= - \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} \cdot \ln \frac{|\boldsymbol{\Sigma}|}{|\mathbf{L}|} \\
&= - \frac{1}{2} \ln \frac{|\boldsymbol{\Sigma}|}{|\mathbf{L}|} = \frac{1}{2} \ln \frac{|\mathbf{L}|}{|\boldsymbol{\Sigma}|}
\end{aligned}$$

となり、第二項については、

$$\begin{aligned}
&\frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x} - \mathbf{m})^T \mathbf{L}^{-1} (\mathbf{x} - \mathbf{m}) d\mathbf{x} \\
&= \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x}^T \mathbf{L}^{-1} - \mathbf{m}^T \mathbf{L}^{-1}) (\mathbf{x} - \mathbf{m}) d\mathbf{x} \\
&= \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x}^T \mathbf{L}^{-1} \mathbf{x} - \mathbf{x}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \mathbf{x} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m}) d\mathbf{x}
\end{aligned}$$

と展開でき、ここで、ベクトルと行列の二次形式の性質より、 $\mathbf{x}^T \mathbf{L}^{-1} \mathbf{x} = \text{Tr}[\mathbf{x}^T \mathbf{L}^{-1} \mathbf{x}]$  が成り立ち、さらにトレースの循環性より、 $\text{Tr}[\mathbf{x}^T \mathbf{L}^{-1} \mathbf{x}] = \text{Tr}[\mathbf{L}^{-1} \mathbf{x} \mathbf{x}^T]$  が成り立つことから、上記の式は、

$$= \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \{ \text{Tr}[\mathbf{L}^{-1} \mathbf{x} \mathbf{x}^T] - \mathbf{x}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \mathbf{x} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} \} d\mathbf{x}$$

と書き直せる。これに対し、式 ( 2.57 ), ( 2.59 ), ( 2.62 ) を適用すると、

$$= \frac{1}{2} ( \text{Tr}[\mathbf{L}^{-1} ( \boldsymbol{\mu} \boldsymbol{\mu}^T + \boldsymbol{\Sigma} )] - \boldsymbol{\mu}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \boldsymbol{\mu} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} )$$

となり、第三項についても同様に、

$$\begin{aligned}
& -\frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) d\mathbf{x} \\
& = \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1}) (\mathbf{x} - \boldsymbol{\mu}) d\mathbf{x} \\
& = \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) d\mathbf{x}
\end{aligned}$$

と展開でき、ここで、ベクトルと行列の二次形式の性質より、 $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} = \text{Tr}[\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}]$  が成り立ち、さらにトレースの循環性より、 $\text{Tr}[\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}] = \text{Tr}[\boldsymbol{\Sigma}^{-1} \mathbf{x} \mathbf{x}^T]$  が成り立つことから、上記の式は、

$$= \frac{1}{2} \int N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) (\text{Tr}[\boldsymbol{\Sigma}^{-1} \mathbf{x} \mathbf{x}^T] - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) d\mathbf{x}$$

と書き直せる。これに対し、式 ( 2.57 ), ( 2.59 ), ( 2.62 ) を適用すると、

$$\begin{aligned}
& = \frac{1}{2} (\text{Tr}[\boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} \boldsymbol{\mu}^T + \boldsymbol{\Sigma})] - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \\
& = \frac{1}{2} (\text{Tr}[\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \boldsymbol{\mu}^T + \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}] - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \\
& = \frac{1}{2} (\text{Tr}[\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \boldsymbol{\mu}^T] + \text{Tr}[\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}] - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})
\end{aligned}$$

となるので、 $\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} = \mathbf{I}_D$  となることと、トレースの循環性とベクトルと行列の二次形式の性質を再度用いると、

$$\begin{aligned}
& = \frac{1}{2} (\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \text{Tr}[\mathbf{I}_D] - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \\
& = \frac{1}{2} \text{Tr}[\mathbf{I}_D] = -\frac{1}{2} D
\end{aligned}$$

となる。以上より、式 ※ で表される 2 つのガウス分布  $p(\mathbf{x}) = N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  と  $q(\mathbf{x}) = N(\mathbf{x} | \mathbf{m}, \mathbf{L})$  の間のカルバック – ライブラーダイバージェンス  $\text{KL}(p \| q)$  は、

$$\text{KL}(p \| q) = \frac{1}{2} \left( \ln \frac{|\mathbf{L}|}{|\boldsymbol{\Sigma}|} + \text{Tr}[\mathbf{L}^{-1} (\boldsymbol{\mu} \boldsymbol{\mu}^T + \boldsymbol{\Sigma})] - \boldsymbol{\mu}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \boldsymbol{\mu} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} - D \right)$$

と求められる。さらに簡潔にまとめると、

$$\begin{aligned}
& = \frac{1}{2} \left( \ln \frac{|\mathbf{L}|}{|\boldsymbol{\Sigma}|} + \text{Tr}[\mathbf{L}^{-1} \boldsymbol{\mu} \boldsymbol{\mu}^T + \mathbf{L}^{-1} \boldsymbol{\Sigma}] - \boldsymbol{\mu}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \boldsymbol{\mu} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} - D \right) \\
& = \frac{1}{2} \left( \ln \frac{|\mathbf{L}|}{|\boldsymbol{\Sigma}|} + \text{Tr}[\mathbf{L}^{-1} \boldsymbol{\Sigma}] + \text{Tr}[\mathbf{L}^{-1} \boldsymbol{\mu} \boldsymbol{\mu}^T] - \boldsymbol{\mu}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \boldsymbol{\mu} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} - D \right) \\
& = \frac{1}{2} \left( \ln \frac{|\mathbf{L}|}{|\boldsymbol{\Sigma}|} + \text{Tr}[\mathbf{L}^{-1} \boldsymbol{\Sigma}] + \text{Tr}[\boldsymbol{\mu}^T \mathbf{L}^{-1} \boldsymbol{\mu}] - \boldsymbol{\mu}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \boldsymbol{\mu} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} - D \right)
\end{aligned}$$

$$= \frac{1}{2} \left( \ln \frac{|\mathbf{L}|}{|\mathbf{\Sigma}|} + \text{Tr}[\mathbf{L}^{-1}\mathbf{\Sigma}] + \boldsymbol{\mu}^T \mathbf{L}^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^T \mathbf{L}^{-1} \mathbf{m} - \mathbf{m}^T \mathbf{L}^{-1} \boldsymbol{\mu} + \mathbf{m}^T \mathbf{L}^{-1} \mathbf{m} - D \right)$$

$$= \frac{1}{2} \left\{ \ln \frac{|\mathbf{L}|}{|\mathbf{\Sigma}|} + \text{Tr}[\mathbf{L}^{-1}\mathbf{\Sigma}] + (\boldsymbol{\mu}^T - \mathbf{m}^T) \mathbf{L}^{-1} (\boldsymbol{\mu} - \mathbf{m}) - D \right\}$$

とまとめられる。

#### [ トレース ( trace ) ]

線形代数において、基底変換を行うと、行列の具体形は変わってしまう。しかし、トレースや行列式は、基底の取り方に依らないので、座標表示の取り方に依らない不変量として非常に重要となる。また、トレースは簡単な形をしているので、計算が楽になる。

$$\circ \text{Tr}(\mathbf{A}) = a_{11} + a_{22} + \cdots + a_{NN}$$

⇒  $N \times N$  正方行列  $\mathbf{A}$  の主対角成分の和で求められる。

#### [ トレースの性質 ]

$$\circ \text{連結性} \cdots \text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B})$$

$$\circ \text{Tr}(k\mathbf{A}) = k \text{Tr}(\mathbf{A})$$

$$\circ \text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$$

$$\circ \det(\mathbf{P}) \neq 0 \text{ のとき、} \text{Tr}(\mathbf{P}^{-1}\mathbf{AP}) = \text{Tr}(\mathbf{A})$$

$$\circ \mathbf{A}^2 - \text{Tr}(\mathbf{A})\mathbf{A} + \det(\mathbf{A})\mathbf{E} = 0$$

$$\circ \text{Tr}(\mathbf{A}^2) - (\text{Tr}(\mathbf{A}))^2 + 2\det(\mathbf{A}) = 0$$

$$\circ \text{循環性} \cdots \text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{CAB}) = \text{Tr}(\mathbf{BCA})$$

⇒ この循環性は、任意の数の行列に対しても拡張される。

(ただし、 $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ ,  $\mathbf{P}$  は  $N$  次の正方行列とする。)

$$\circ \text{Tr}(\mathbf{ab}^T) = \text{Tr}\{(\mathbf{ab}^T)^T\} = \text{Tr}(\mathbf{ba}^T)$$

⇒ トレース内で転置をとっても等しくなる。

$$\circ \text{ベクトルと行列の二次形式の性質} \cdots \mathbf{a}^T \mathbf{\Sigma} \mathbf{a} = \text{Tr}(\mathbf{a}^T \mathbf{\Sigma} \mathbf{a})$$

(ただし、 $\mathbf{a}$ ,  $\mathbf{b}$  は  $N$  次のベクトル、 $\mathbf{\Sigma}$  は  $N$  次の正方行列とする。)