

## Ch.3 Bayes推論の基礎

### §3.1 確率推論

#### 3.1.1 確率密度関数と確率質量関数

$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^M$  の実数値関数  $p(x)$  が次の二つの条件をみたすとき,

$p(x)$  を **確率密度関数** (probability density function: pdf) という.

(1)  $p(x) \geq 0$ .

(2)  $\int p(x) dx = 1$ .

各要素が離散値の  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{Z}^M$  に対して,

実数値関数  $p(x)$  が次の二つの条件をみたすとき,  $p(x)$  を

**確率質量関数** (probability mass function: pmf) という.

(1)  $p(x) \geq 0$ .

(2)  $\sum_x p(x) = 1$ .

pdf や pmf で決まる  $x$  の分布を **確率分布** (probabilistic distribution),

あるいは **確率モデル** (probabilistic model) という.

#### 3.1.2 条件付き分布と周辺分布

二つの変数  $x, y$  に対する確率分布  $p(x, y)$  を **同時分布** (joint distribution)

という.  $p(y) = \int p(x, y) dx$

のようにして一方の変数を積分で除去する操作を **周辺化** (marginalization)

という.  $p(y)$  を  $y$  の **周辺分布** (marginal distribution) という.

- 同時分布  $p(x, y)$  で  $y$  に対して特定の値が決められたときの  $x$  の確率分布を **条件付き分布** (Conditional distribution) といい.

次に定義する:

$$p(x|y) = \frac{p(x, y)}{p(y)}$$

- $y$  は  $x$  の分布  $p(x|y)$  の特性を決めるパラメータのようなものと解釈できる.

- $p(x|y)$  は確かに pdf になっている.

pf.  $p(x|y) \geq 0$  は明らか.

$$\int p(x|y) dx = \frac{1}{p(y)} \int p(x, y) dx = \frac{1}{p(y)} p(y) = 1. \quad \blacksquare$$

- 同時分布  $p(x, y)$

$$p(x, y) = p(x)p(y)$$

とみたとき,  $x$  と  $y$  は **独立** (independent) という.

- ある同時分布が与えられたときに, それから興味の対象となる条件付き分布や周辺分布を算出すること (Bayes) **推論** という.

### 3.1.3 期待値

- 分布  $p(x)$  に対して, 関数  $f(x)$  の **期待値** (expectation) は

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

と計算される.  $\hookrightarrow \mathbb{E}_p[f(x)]$  とか  $\mathbb{E}[f(x)]$  とか書くこともある.

二つの確率分布  $p(x)$  と  $q(x)$  に対し,  $q$  の  $p$  に対する

Kullback-Leibler ダイバージェンス (Kullback-Leibler divergence) と

$$D_{KL}[q(x) \parallel p(x)] := \int q(x) \log \frac{q(x)}{p(x)} dx \\ = \mathbb{E}_q[\log q(x)] - \mathbb{E}_q[\log p(x)]$$

で定義する.

「 $x$  が  $p$  よりも  $q$  から出たと考えられる割合」. より正確には  
「帰無仮説  $p(x)$  を棄却しているかどうかを  
を調べるための統計量」

•  $\frac{q(x)}{p(x)}$ : 尤度比 (likelihood ratio).

KL-divergence は  $q$  を真の分布とみたときの対数尤度比  $\frac{q(x)}{p(x)}$  の平均.

→ 「KL-divergence が大きい」とは「平均的に分布  $q$  の方が  $p$  よりも

尤もらしい」ということ.  $p$  は  $q$  の分布から「かなり離れたところ」.

• もう少し深入りしてみる.

•  $H(q) := -\mathbb{E}_q[\log q(x)]$  は分布  $q$  のエントロピー (entropy).

分布  $q$  の平均の情報量であり,  $q$  に従う事象の予測の難しさを表す.

ex)  $x \in \{0, 1\}$ .  $p(x) = \begin{cases} 0.9 & (x=1) \\ 0.1 & (x=0) \end{cases}$  ←  $p$  に従うときは  
「 $x=1$  になる」.

$q(x) = 0.5$  ←  $x=0$  と  $x=1$  はおきりしない.

$$H(p) = -0.9 \log 0.9 - 0.1 \log 0.1 \approx 0.47$$

$$H(q) = -0.5 \log 0.5 - 0.5 \log 0.5 \approx 0.72$$

$H(q) > H(p)$  なること

←  $q$  の方が予測の難しい分布!

•  $H(q, p) := -\mathbb{E}_q[\log p(x)]$  は分布  $p$  の  $q$  に対する **交差エントロピー** (Cross entropy)

真の分布  $p$  とみたときの負の対数尤度  $-\log p(x)$  の平均.

→ (小さいほど) 分布  $p$  が尤もらしい, i.e.,  $p$  が  $q$  に似ている.

•  $D_{KL}[q(x) \| p(x)] = H(q, p) - H(q).$

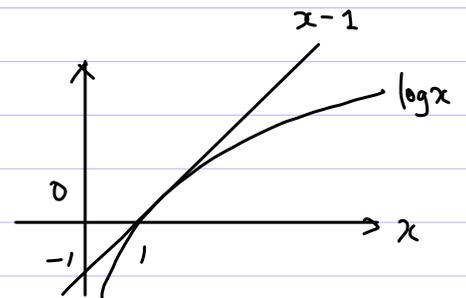
KL-divergence は, 「 $p$  の  $q$  からの解離度」に「 $q$  の複雑さ」を加味した値になっている.

Th. 任意の  $p, q$  に対して  $D_{KL}[q(x) \| p(x)] \geq 0$ . (**Gibbsの不等式**)  
(等号成立は  $q(x) = p(x), \forall x$  のとき)

pf. 以下の不等式を利用する.

$$\log x \leq x - 1 \quad (x > 0)$$

(等号成立は  $x=1$  のとき.)



$$\begin{aligned} D_{KL}[q(x) \| p(x)] &= \int q(x) \log \frac{q(x)}{p(x)} dx \\ &= - \int q(x) \log \frac{p(x)}{q(x)} dx \\ &\geq - \int q(x) \left( \frac{p(x)}{q(x)} - 1 \right) dx \\ &= \underbrace{- \int p(x) dx}_{=1} + \underbrace{\int q(x) dx}_{=1} = 0. \end{aligned}$$

等号成立は  $\frac{p(x)}{q(x)} = 1 \Leftrightarrow p(x) = q(x)$  のとき.

cf. Jensen の不等式を用いて示せる.



• KL-divergenceは非対称的:  $D_{KL}[p||q] \neq D_{KL}[q||p]$ .

(定義式から分かる)

• 三角不等式も満たさない. 実は「距離の2乗」のような性質をもつ.  
詳しくは情報幾何の本を見よ

→ KL-divergenceは「距離」ではない.

cf. 後々に「真の分布  $p$  に近似分布  $q$  を近づけたために

$D_{KL}[q||p]$  を最小化する」といった操作が出る.

上では「 $q$  を真の分布と考えて  $p$  から  $q$  へどれだけの距離があるか」と

見ていたのでは「違和感」がある.  $D_{KL}[p||q]$  を最小化するべきでは?

•  $D_{KL}[p||q]$  は forward KL-divergence,

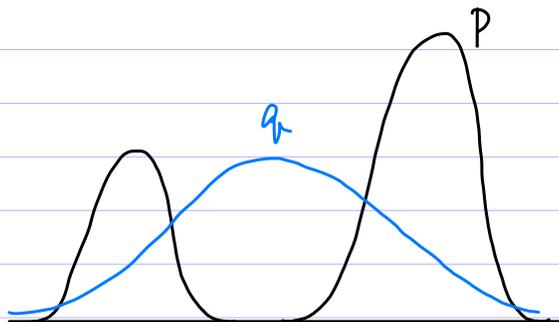
$D_{KL}[q||p]$  は reverse KL-divergence という.

• forward divergence  $D_{KL}[p||q]$  では,  $q(x) = 0$  となる  $x$  で 「絶対連続性」

被積分関数が発散する.  $D_{KL}[p||q]$  の最小化では  $q(x) = 0$  となる

$x$  は対して  $p(x) = 0$  となるべきではない (対偶をとれば,  $p(x) > 0$

ならば  $q(x) > 0$  とする).  $p$  の台 (正值となる区間) を覆うように  $q$  が決定.



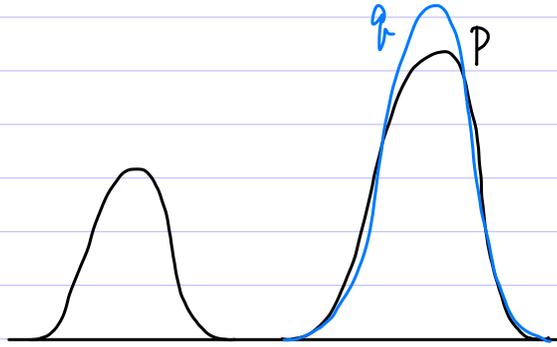
$q$ : Gauss分布

$p$  に近似.

一方,  $D_{KL}[q||p]$  では,  $p(x)=0$  となる  $x$  を 棄却するのを,

最小化すると  $p(x)=0$  となる  $x$  を  $q(x)=0$  とする

( $q(x)>0$  かつ  $p(x)>0$ ).  $p$  の台の一部にフィットするように  $q$  を決定



$q$ : Gauss分布

$p$  に近似.

とすると最小化する  $p$  は, やりたことによる.

$D_{KL}[q||p]$  の近似計算は,  $q$  からのサンプリングが必要.

→  $q$  は簡単な分布でないとい計算が必要.

### 3.1.4 変数変換.

$f: \mathbb{R}^M \rightarrow \mathbb{R}^M$ : bijective (逆写像  $f^{-1}$  がある)

$y = f(x)$  と変数変換する.

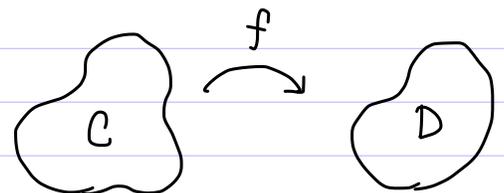
$x$  についての pdf  $p_x(x)$  とする.  $y$  の pdf  $p_y(y)$  はどう書けるか?

$D \subseteq \mathbb{R}^M$  とする.  $C := \{x \in \mathbb{R}^M \mid f(x) \in D\} = f^{-1}(D)$  とする.

$$P(y \in D) = P(f(x) \in D)$$

$$= P(x \in f^{-1}(D))$$

$$= P(x \in C)$$



$$P(y \in D) = \int_D p_y(y) dy.$$

$$P(x \in C) = \int_C p_x(x) dx \quad \text{重積分の変数変換.}$$

$$= \int_D p_x(f^{-1}(y)) \left| \det \frac{\partial x}{\partial y} \right| dy. \quad (*)$$

$\therefore \frac{\partial x}{\partial y}$  は  $f^{-1}$  の **Jacobi 行列** である.

$$\frac{\partial x}{\partial y} := \begin{pmatrix} \frac{\partial x_1}{\partial y_1} & \cdots & \frac{\partial x_1}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial x_n}{\partial y_1} & \cdots & \frac{\partial x_n}{\partial y_n} \end{pmatrix} = \left( \frac{\partial x_i}{\partial y_j} \right)_{ij}.$$

$\det \frac{\partial x}{\partial y}$  は **Jacobian** と呼ばれる。  
 変換の「体積拡大率」を表す。

(\*) の式より、変数変換の公式

$$p_y(y) = p_x(f^{-1}(y)) \left| \det \frac{\partial x}{\partial y} \right|$$

$p_y(y) dy = p_x(x) dx$   
 と覚えておくと思い出しやすいかも

を得る。

ex) Gauss分布に従う r.v.  $x$  を  $\tanh$  により変換して  $y = \tanh(x)$  の分布の pdf を求める。

$$x \text{ に対応する pdf: } p_x(x) = \mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right).$$

$$\frac{dy}{dx} = \frac{d}{dx} \left( \frac{e^x - e^{-x}}{e^x + e^{-x}} \right) = \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} = 1 - \tanh^2(x).$$

$$= 1 - y^2. \quad (> 0)$$

$$\therefore \frac{dx}{dy} = \left( \frac{dy}{dx} \right)^{-1} = \frac{1}{1 - y^2}$$

$$\therefore p_y(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(\tanh^{-1}(y) - \mu)^2\right) \frac{1}{1 - y^2}.$$

$$y = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \Leftrightarrow (e^x + e^{-x})y = e^x - e^{-x}$$

$$\Leftrightarrow (e^{2x} + 1)y = e^{2x} - 1$$

$$\Leftrightarrow e^{2x}(y-1) = -y-1$$

$$\Leftrightarrow e^{2x} = \frac{1+y}{1-y}$$

$$\Leftrightarrow 2x = \log \frac{1+y}{1-y}$$

$$\Leftrightarrow x = \frac{1}{2} \log \frac{1+y}{1-y}.$$

$$\therefore \tanh^{-1}(y) = \frac{1}{2} \log \frac{1+y}{1-y}.$$

$$\therefore p_y(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \left(\frac{1}{2} \log \frac{1+y}{1-y} - \mu\right)^2\right) \frac{1}{1-y^2}.$$

### 3.1.5 グラフィカルモデル.

#### ・ グラフィカルモデル (graphical model):

確率モデルに存在する変数の関係性を有向グラフで表したものを.

**DAG** (directed acyclic graph: 非閉路的有向グラフ) による表現を

考える.

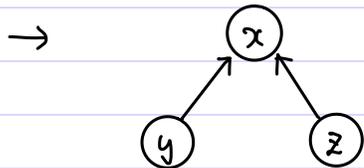
・ 頂点は変数に対応. 辺は変数間の依存関係を表す.

・ ある変数  $x_i$  与えらるると  $y$  の分布が定まる ( $p(y|x_i)$ :  $y$  は  $x_i$  による

条件付けらるる) とし, 頂点  $x_i$  から  $y$  へ向かう有向辺を張る.

ex) 3変数  $x, y, z$  からなる確率モデル

$$p(x, y, z) = p(x|y, z) p(y) p(z)$$

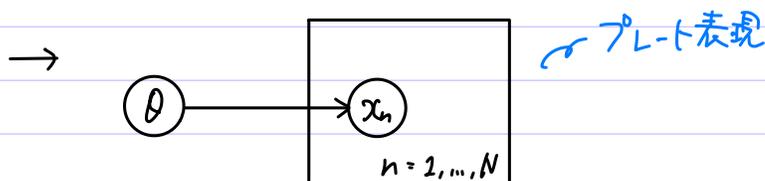


ex) パラメータ  $\theta$  に依存して  $N$  個の変数  $\mathcal{X} = \{x_1, \dots, x_N\}$  が発生する

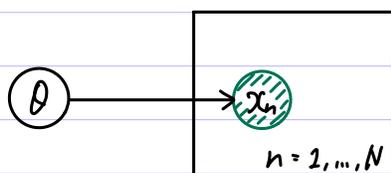
確率モデル

$$p(\mathcal{X}, \theta) = \underbrace{p(\mathcal{X}|\theta)}_{\text{尤度関数 (likelihood function)}} \underbrace{p(\theta)}_{\text{事前分布 (prior distribution)}}$$

$= \left( \prod_{n=1}^N p(x_n|\theta) \right) p(\theta)$



$\mathcal{X}$ : 観測データ のときは上の区間



などと塗ることで観測されていることを表すことがある。



$$= \mu \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) dx + \int_{-\infty}^{\infty} (x-\mu) \mathcal{N}(x|\mu, \sigma^2) dx$$

= 1 (∵ pdfの全領域での積分)

$$= \mu + \int_{-\infty}^{\infty} (x-\mu) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx$$

$$= \mu + \sigma^2 \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \left( \frac{\partial}{\partial \mu} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \right) dx$$

$$= \mu + \sigma^2 \frac{\partial}{\partial \mu} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx$$

= 1 (∵ 被積分関数は  $\mathcal{N}(x|\mu, \sigma^2)$ )

微分と積分の  
順序交換.

$$= \mu + \sigma^2 \left( \frac{\partial}{\partial \mu} 1 \right)$$

$$= \mu. \quad \blacksquare$$

$$\cdot \mathbb{V}_{\mathcal{N}(\mu, \sigma^2)}[x] = \sigma^2.$$

pf. (ゴリ押し)

$$\mathbb{V}_{\mathcal{N}(\mu, \sigma^2)}[x]$$

$$= \int_{-\infty}^{\infty} (x-\mu)^2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx \quad \downarrow t = \frac{x-\mu}{\sqrt{2\sigma^2}}$$

$$= \frac{2\sigma^2}{\sqrt{\pi}} \int_{-\infty}^{\infty} t^2 e^{-t^2} dt \quad \downarrow \text{部分積分}$$

$$= \frac{2\sigma^2}{\sqrt{\pi}} \left( \underbrace{\left[ t \cdot \left(-\frac{1}{2} e^{-t^2}\right) \right]_{-\infty}^{\infty}}_{=0} + \frac{1}{2} \int_{-\infty}^{\infty} \underbrace{e^{-t^2}}_{=\sqrt{\pi}} dt \right)$$

$$= \frac{2\sigma^2}{\sqrt{\pi}} \cdot \frac{1}{2} \sqrt{\pi}$$

$$= \sigma^2. \quad \blacksquare$$

(パラメータによる微分を用いる)

以下では 精度パラメータ  $\lambda := \frac{1}{\sigma^2}$  による微分を利用する.

$$\mathcal{N}(x|\mu, \lambda) = \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda}{2}(x-\mu)^2\right).$$

$$\begin{aligned}
& V_{\mathcal{N}(\mu, \sigma^2)}[x] \\
&= \int_{-\infty}^{\infty} (x-\mu)^2 \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda}{2}(x-\mu)^2\right) dx \\
&= \sqrt{\frac{\lambda}{2\pi}} \int_{-\infty}^{\infty} (x-\mu)^2 \exp\left(-\frac{\lambda}{2}(x-\mu)^2\right) dx \\
&= \sqrt{\frac{\lambda}{2\pi}} \cdot (-2) \int_{-\infty}^{\infty} \left(\frac{\partial}{\partial \lambda} \exp\left(-\frac{\lambda}{2}(x-\mu)^2\right)\right) dx \\
&= -2 \sqrt{\frac{\lambda}{2\pi}} \frac{\partial}{\partial \lambda} \int_{-\infty}^{\infty} \exp\left(-\frac{\lambda}{2}(x-\mu)^2\right) dx \quad \left\} \begin{array}{l} \text{微分と積分の順序交換} \\ \text{Gauss 積分} \end{array} \right. \\
&= -2 \sqrt{\frac{\lambda}{2\pi}} \frac{\partial}{\partial \lambda} \sqrt{\frac{2\pi}{\lambda}} \\
&= -2\sqrt{\lambda} \left(-\frac{1}{2} \frac{1}{(\sqrt{\lambda})^3}\right) \\
&= \frac{1}{\lambda} \\
&= \sigma^2. \quad \blacksquare
\end{aligned}$$

•  $M$ -次元 Gauss 分布.

$$\text{pdf: } \mathcal{N}(x | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^M \det \Sigma}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

$\mu \in \mathbb{R}^M$ : 平均ベクトル.

$\Sigma \in M_M(\mathbb{R})$ : 正定値 ( $\Sigma > 0$  と書く) 共分散行列

$$\cdot \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[x] = \mu.$$

$$\text{pf. } \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[x]$$

$$= \int_{\mathbb{R}^M} x \mathcal{N}(x | \mu, \Sigma) dx$$

$$= \mu \underbrace{\int_{\mathbb{R}^M} \mathcal{N}(x | \mu, \Sigma) dx}_{=1} + \int_{\mathbb{R}^M} (x-\mu) \mathcal{N}(x | \mu, \Sigma) dx$$

= 1

$$= \mu + \int_{\mathbb{R}^M} (x - \mu) \frac{1}{\sqrt{(2\pi)^M \det \Sigma}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) dx.$$

$\Sigma^{-1}$ : 対称  $\Sigma$ : 対称  $\Sigma^{-1}$ : 対称.  $\leftarrow \Sigma = \Sigma^T$  より  $(\Sigma^{-1})^T = (\Sigma^T)^{-1} = \Sigma^{-1}$ .

微分公式  $\frac{\partial}{\partial x}(x^T A x) = 2Ax$  ( $A$ : 対称) を使えば,

$$\begin{aligned} & \frac{\partial}{\partial \mu} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \quad \downarrow \text{合成関数の微分.} \\ &= \frac{\partial}{\partial \mu} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \frac{\partial}{\partial \mu} \left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \quad \downarrow \text{合成関数の微分.} \\ &= \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \left(-\frac{1}{2} \cdot 2 \Sigma^{-1}(x - \mu) \frac{\partial}{\partial \mu}(x - \mu)\right) \\ &= \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \left(-\Sigma^{-1}(x - \mu) (-I_M)\right) \quad \begin{array}{l} \uparrow \\ \text{ベクトルのベクトルの微分.} \\ \rightarrow \text{Jacobi行列.} \end{array} \\ &= \Sigma^{-1}(x - \mu) \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right). \end{aligned}$$

よって,

$$\begin{aligned} & \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}[x] \\ &= \mu + \Sigma \int_{\mathbb{R}^M} \frac{1}{\sqrt{(2\pi)^M \det \Sigma}} \left( \frac{\partial}{\partial \mu} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \right) dx \\ &= \mu + \Sigma \frac{\partial}{\partial \mu} \int_{\mathbb{R}^M} \frac{1}{\sqrt{(2\pi)^M \det \Sigma}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) dx \\ & \quad = 1 \quad (\because \text{被積分関数は pdf}) \\ &= \mu + \Sigma \left( \frac{\partial}{\partial \mu} 1 \right) \\ &= \mu + \Sigma 0 \\ &= \mu. \quad \square \end{aligned}$$

$$\cdot \mathcal{N}(\mu, \Sigma)(x) = \Sigma.$$

pf. まず, 以下を Fact とし使おう (Cholesky分解)

[Fact]  $A > 0 \Leftrightarrow \exists L: \text{下三角行列, 対角項が正 s.t.}$

$$A = LL^T.$$

$\Sigma > 0$  かつ  $\Sigma^{-1} > 0$ .  $\leftarrow \Sigma$  の固有値  $\lambda_i > 0$  ( $i=1, \dots, M$ ) かつ  $\Sigma^{-1}$  の固有値  $\lambda_i^{-1} > 0$  ( $i=1, \dots, M$ ) かつ.

$\Sigma^{-1} = LL^T$  と Cholesky分解すれば  $\Sigma^{-1}$  は  $L$  の逆行列と見做す.

$$\begin{aligned} \mathcal{N}(x | \mu, L) &= \frac{\det L}{\sqrt{(2\pi)^M}} \exp\left(-\frac{1}{2} (x-\mu)^T L L^T (x-\mu)\right) \\ &= \frac{\det L}{\sqrt{(2\pi)^M}} \exp\left(-\frac{1}{2} \text{tr}\left((x-\mu)^T L L^T (x-\mu)\right)\right) \\ &= \frac{\det L}{\sqrt{(2\pi)^M}} \exp\left(-\frac{1}{2} \text{tr}\left(L^T (x-\mu)(x-\mu)^T L\right)\right). \end{aligned}$$

$\text{tr}$  の性質  
 $\text{tr}(ABC) = \text{tr}(CAB)$ .

[Claim]  $A \in M_n(\mathbb{R})$ : 対称 とし  $\frac{\partial}{\partial L} \text{tr}(L^T A L) = 2AL$ .

pf.  $L = (l_{ij}), A = (a_{ij})$  とし

$$\begin{aligned} \frac{\partial}{\partial l_{ij}} \text{tr}(L^T A L) &= \frac{\partial}{\partial l_{ij}} \sum_{m=1}^M (L^T A L)_{mm} \\ &= \frac{\partial}{\partial l_{ij}} \sum_{m=1}^M \sum_{s=1}^M \sum_{t=1}^M l_{sm} a_{st} l_{tm} \\ &= \sum_{m=1}^M \sum_{s=1}^M \sum_{t=1}^M \left( \frac{\partial l_{sm}}{\partial l_{ij}} a_{st} l_{tm} + l_{sm} a_{st} \frac{\partial l_{tm}}{\partial l_{ij}} \right) \\ &= \sum_{t=1}^M a_{itt} l_{tj} + \sum_{s=2}^M l_{sj} a_{si} \\ &= 2 \sum_{k=1}^M a_{ik} l_{kj} \\ &= 2(AL)_{ij}. \end{aligned}$$

以上の Claim を利用すると,

$$\begin{aligned}
 & \frac{\partial}{\partial L} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) \quad \downarrow \text{合成関数の微分.} \\
 &= \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) \frac{\partial}{\partial L} \left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) \\
 &= -\frac{1}{2} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) \frac{\partial}{\partial L} \left(\text{tr}\left(L^T \underbrace{(\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^T}_{\text{対称行列}} L\right)\right) \\
 &= -\cancel{\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) \cdot \cancel{2} (\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^T L \\
 &= -(\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^T \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) L.
 \end{aligned}$$

よって,

$$\begin{aligned}
 & \mathbb{V}_{N(\boldsymbol{\mu}, L)}[\mathbf{x}] \\
 &= \int_{\mathbb{R}^n} (\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^T \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, L) d\mathbf{x} \\
 &= \int_{\mathbb{R}^n} (\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^T \frac{\det L}{\sqrt{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) d\mathbf{x} \\
 &= \frac{\det L}{\sqrt{(2\pi)^n}} \left( \int_{\mathbb{R}^n} -\frac{\partial}{\partial L} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) d\mathbf{x} \right) L^{-1} \quad \downarrow \text{微分と積分の順序交換} \\
 &= -\frac{\det L}{\sqrt{(2\pi)^n}} \left( \frac{\partial}{\partial L} \int_{\mathbb{R}^n} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T L L^T (\mathbf{x}-\boldsymbol{\mu})\right) d\mathbf{x} \right) L^{-1} \\
 &= -\frac{\det L}{\sqrt{(2\pi)^n}} \left( \frac{\partial}{\partial L} \frac{\sqrt{(2\pi)^n}}{\det L} \int_{\mathbb{R}} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, L) d\mathbf{x} \right) L^{-1} \\
 &= -\det L \left( \frac{\partial}{\partial L} (\det L)^{-1} \right) L^{-1} \quad = 1 \\
 &= \det L \cdot \frac{1}{(\det L)^2} \left( \frac{\partial}{\partial L} \det L \right) L^{-1} \\
 &= (\det L)^{-1} \left( \frac{\partial}{\partial L} \det L \right) L^{-1}.
 \end{aligned}$$

[Claim]  $\frac{\partial}{\partial L} \det L = (\det L) (L^{-1})^T$ .

pf.  $L_{ij}$ :  $L$  の第  $i$  行第  $j$  列を除いてできる小行列.

$\Delta_{ij} := (-1)^{i+j} \det L_{ij}$  :  $(i, j)$ -余因子.  $L_{ij}$  に依る  $L$

第  $i$  行についての Laplace 展開:  $\det L = \sum_{t=1}^M l_{it} \Delta_{it}$

$\Sigma$  用いて

$$\frac{\partial}{\partial l_{ij}} \det L = \sum_{t=1}^M \frac{\partial l_{it}}{\partial l_{ij}} \Delta_{it} = \sum_{t=1}^M \delta_{it} \delta_{tj} \Delta_{it} = \Delta_{ij}$$

$L$  の余因子行列  $\text{Cof}(L) := (\Delta_{ij})$  に対しては,

$$L^{-1} = \frac{1}{\det L} (\text{Cof}(L))^T$$

が成立するの2,  $\frac{\partial}{\partial l_{ij}} \det L = (\det L) ((L^{-1})^T)_{ij}$ .  $\blacksquare$

$\therefore V_{N(\mu, \Sigma)}[x]$

$$= (\det L)^{-1} (\det L) (L^{-1})^T L^{-1}$$

$$= (L^T)^{-1} L^{-1} \quad \downarrow (AB)^{-1} = B^{-1}A^{-1}$$

$$= (LL^T)^{-1}$$

$$= (\Sigma^{-1})^{-1}$$

$$= \Sigma. \quad \blacksquare$$

cf. 直接

$$V_{N(\mu, \Sigma)}[x] = \int_{\mathbb{R}^M} (x-\mu)(x-\mu)^T N(x|\mu, \Sigma) dx$$

$\Sigma$  計算するのより.  $\Sigma^{-1}$ : 対称な2,  $\exists U$ : 直交行列 s.t.

$$\Sigma^{-1} = U \text{diag}(\lambda_1^{-1}, \dots, \lambda_M^{-1}) U^T. \quad \text{これを} \Sigma \text{用いて}$$

$$\Sigma^{-1} = U \text{diag}(\lambda_1^{-\frac{1}{2}}, \dots, \lambda_M^{-\frac{1}{2}}) U^T \cdot U \text{diag}(\lambda_1^{\frac{1}{2}}, \dots, \lambda_M^{\frac{1}{2}}) U^T =: \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}}$$

$(\Sigma^{-\frac{1}{2}} > 0 \text{ 1: } \Sigma \text{ は正定}) \quad \mathbf{z} = \mathbf{z}'$ ,  $\mathbf{Z} := \Sigma^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu})$  とおくと,

$$\mathbf{x} - \boldsymbol{\mu} = \Sigma^{\frac{1}{2}} \mathbf{Z}, \quad (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T = \Sigma^{\frac{1}{2}} \mathbf{Z} \mathbf{Z}^T \Sigma^{\frac{1}{2}}$$

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\Sigma^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}))^T \Sigma^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}) = \mathbf{Z}^T \mathbf{Z}$$

変数変換の Jacobian は  $\det(AB) = \det A \cdot \det B$  などを利用

$$\det \frac{\partial \mathbf{x}}{\partial \mathbf{z}} = \det(\Sigma^{\frac{1}{2}}) = (\det \Sigma)^{\frac{1}{2}}$$

$$\therefore \mathbb{V}_{\mathcal{N}(\boldsymbol{\mu}, \Sigma)}[\mathbf{x}]$$

$$= \Sigma^{\frac{1}{2}} \left( \int_{\mathbb{R}^M} \mathbf{z} \mathbf{z}^T \frac{1}{\sqrt{(2\pi)^M (\det \Sigma)}} \exp\left(-\frac{1}{2} \mathbf{z}^T \mathbf{z}\right) \cdot (\det \Sigma)^{\frac{1}{2}} d\mathbf{z} \right) \Sigma^{\frac{1}{2}}$$

$$= \frac{1}{\sqrt{(2\pi)^M}} \Sigma^{\frac{1}{2}} \left( \int_{\mathbb{R}^M} \mathbf{z} \mathbf{z}^T \exp\left(-\frac{1}{2} \mathbf{z}^T \mathbf{z}\right) d\mathbf{z} \right) \Sigma^{\frac{1}{2}}$$

$$\int_{\mathbb{R}^M} \mathbf{z} \mathbf{z}^T \exp\left(-\frac{1}{2} \mathbf{z}^T \mathbf{z}\right) d\mathbf{z} = \sqrt{(2\pi)^M} \mathbf{I}_M \text{ とおこう! (*exercise*)}$$

$$\mathbb{V}_{\mathcal{N}(\boldsymbol{\mu}, \Sigma)}[\mathbf{x}] = \Sigma$$

$\mathbf{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_M \end{pmatrix}$  とおいて頑張って計算する。

◦ Bernoulli 分布. 2 値  $\Sigma$  と変数  $x \in \{0, 1\}$  を生成するための分布.

$$\text{pmf: } \text{Bern}(x | \mu) = \mu^x (1-\mu)^{1-x}, \quad \mu \in (0, 1) : \text{"成り立確率"を表す.}$$

$$\cdot \mathbb{E}_{\text{Bern}(\mu)}[x] = \mu$$

$$\text{pf. } \mathbb{E}_{\text{Bern}(\mu)}[x] = \sum_{x=0}^1 x \mu^x (1-\mu)^{1-x} = 1 \cdot \mu^1 (1-\mu)^{1-1} = \mu. \quad \square$$

$$\cdot \mathbb{V}_{\text{Bern}(\mu)}[x] = \mu(1-\mu).$$

$$\text{pf. } \mathbb{E}_{\text{Bern}(\mu)}[x(x-1)] = 0.$$

$$\mathbb{E}_{\text{Bern}(\mu)}[x^2] - \mathbb{E}_{\text{Bern}(\mu)}[x]$$

$$\mathbb{V}_{\text{Bern}(\mu)}[x] = \mathbb{E}_{\text{Bern}(\mu)}[x(x-1)] - \mathbb{E}_{\text{Bern}(\mu)}[x]^2 + \mathbb{E}_{\text{Bern}(\mu)}[x] = \mu(1-\mu). \quad \square$$

- **カテゴリ分布**. Bernoulli分布をD値に拡張したものである.

$$S \in \{0,1\}^D \text{ かつ } \sum_{d=1}^D S_d = 1 \text{ とする one-hot ベクトルを生成する.}$$

$$\text{pmf: } \text{Cat}(S|\pi) = \prod_{d=1}^D \pi_d^{S_d}, \quad \pi = \begin{pmatrix} \pi_1 \\ \vdots \\ \pi_D \end{pmatrix} \in (0,1)^D \text{ かつ } \sum_{d=1}^D \pi_d = 1.$$

$$\cdot \mathbb{E}_{\text{Cat}(\pi)}[S] = \pi$$

$$\text{pf. } \mathcal{S} := \{e_d \mid d=1, \dots, D\} \quad (e_d = (\delta_{di})_{i=1}^D : \text{標準基底})$$

$$\mathbb{E}_{\text{Cat}(\pi)}[S] = \sum_{S \in \mathcal{S}} S \prod_{d=1}^D \pi_d^{S_d} = \sum_{d=1}^D \pi_d e_d = \pi. \quad \square$$

$$\cdot \text{Var}_{\text{Cat}(\pi)}[S] = \text{diag}(\pi) - \pi\pi^T.$$

$$\text{pf. } \mathbb{E}_{\text{Cat}(\pi)}[SS^T] = \sum_{S \in \mathcal{S}} SS^T \prod_{d=1}^D \pi_d^{S_d} = \sum_{d=1}^D \pi_d e_d e_d^T = \text{diag}(\pi).$$

$$\begin{aligned} \therefore \text{Var}_{\text{Cat}(\pi)}[S] &= \mathbb{E}_{\text{Cat}(\pi)}[SS^T] - \mathbb{E}_{\text{Cat}(\pi)}[S]\mathbb{E}_{\text{Cat}(\pi)}[S]^T \\ &= \text{diag}(\pi) - \pi\pi^T. \quad \square \end{aligned}$$

- **ガンマ分布**.  $\lambda > 0$  を生成する分布.

$$\text{pdf: } \text{Gam}(\lambda|a,b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}. \quad a, b > 0.$$

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt \quad (x > 0) : \text{ガンマ関数}$$

$$\cdot \Gamma(x+1) = x\Gamma(x). \quad \leftarrow x: \text{整数のとき, } \Gamma(x+1) = x!$$

$$\text{pf. } \Gamma(x+1)$$

$$= \int_0^\infty t^x e^{-t} dt$$

$$= \left[ -t^x e^{-t} \right]_0^\infty + x \int_0^\infty t^{x-1} e^{-t} dt$$

$$= x\Gamma(x). \quad \square$$

ガンマ関数は階乗の一般化

$$\cdot \mathbb{E}_{\text{Gam}(a,b)}[\lambda] = \frac{a}{b}$$

pf.  $\mathbb{E}_{\text{Gam}(a,b)}[\lambda]$

$$= \int_0^{\infty} \lambda \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} d\lambda$$

$$= \frac{b^a}{\Gamma(a)} \int_0^{\infty} \lambda^a e^{-b\lambda} d\lambda$$

$$= \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+1)}{b^{a+1}} \int_0^{\infty} \text{Gam}(\lambda | a+1, b) d\lambda \quad \left. \begin{array}{l} \text{=} 1 \\ \downarrow \Gamma(a+1) = a\Gamma(a) \end{array} \right\}$$

$$= \frac{a}{b} \quad \blacksquare$$

$$\cdot \mathbb{V}_{\text{Gam}(a,b)}[\lambda] = \frac{a}{b^2}$$

pf.  $\mathbb{E}_{\text{Gam}(a,b)}[\lambda^2]$

$$= \int_0^{\infty} \lambda^2 \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} d\lambda$$

$$= \frac{b^a}{\Gamma(a)} \int_0^{\infty} \lambda^{a+1} e^{-b\lambda} d\lambda$$

$$= \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+2)}{b^{a+2}}$$

$$= \frac{(a+1)a}{b^2}$$

$$\therefore \mathbb{V}_{\text{Gam}(a,b)}[\lambda] = \mathbb{E}_{\text{Gam}(a,b)}[\lambda^2] - \mathbb{E}_{\text{Gam}(a,b)}[\lambda]^2 = \frac{a}{b^2} \quad \blacksquare$$

### 3.2.2 Gauss分布の計算例.

$$\cdot \mathbf{x} \in \mathbb{R}^D, \quad \mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}, \quad \mathbf{x}_1 \in \mathbb{R}^{D_1}, \quad \mathbf{x}_2 \in \mathbb{R}^{D_2} \quad (D = D_1 + D_2) \text{ と } \exists.$$

$$p(\mathbf{x}) = p(\mathbf{x}_1, \mathbf{x}_2) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ と } \exists. \quad (D\text{-次元 Gauss 分布})$$

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \quad \begin{array}{l} \boldsymbol{\mu}_1 \in \mathbb{R}^{D_1} \\ \boldsymbol{\mu}_2 \in \mathbb{R}^{D_2} \end{array} \quad \boldsymbol{\Sigma} = \left( \begin{array}{c|c} \boldsymbol{\Sigma}_{1,1} & \boldsymbol{\Sigma}_{1,2} \\ \hline \boldsymbol{\Sigma}_{2,1} & \boldsymbol{\Sigma}_{2,2} \end{array} \right) \begin{array}{l} D_1 \\ D_2 \end{array} \quad (\boldsymbol{\Sigma}_{1,2} = \boldsymbol{\Sigma}_{2,1}^T)$$

$\boldsymbol{\Sigma}$ : 対称テンソル.

精度行列 (precision matrix) .  $\Lambda := \Sigma^{-1}$ .

$$\Lambda = \left( \begin{array}{c|c} \Lambda_{1,1} & \Lambda_{1,2} \\ \hline \Lambda_{2,1} & \Lambda_{2,2} \end{array} \right) \begin{array}{l} D_1 \\ D_2 \end{array} \quad (\Lambda_{1,2} = \Lambda_{2,1}^T)$$

$D_1$        $D_2$

Th. (ブロック行列の逆行列)

$$A \in M_n(\mathbb{R}), B \in M_{n,m}(\mathbb{R}), C \in M_{m,n}(\mathbb{R}),$$

$$D \in M_m(\mathbb{R}) : \text{regular}.$$

$$M := A - BD^{-1}C \quad (\text{Schur 補行列}) : \text{regular とする}.$$

$$\Rightarrow \begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} M^{-1} & -M^{-1}BD^{-1} \\ -D^{-1}CM^{-1} & D^{-1} + D^{-1}CM^{-1}BD \end{pmatrix}.$$

pf. 行列  $\begin{pmatrix} A & B \\ C & D \end{pmatrix}$  の基本変形 E 行う

以下,  $\rightarrow$  で基本変形 E 表すものとする.

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} \begin{array}{l} \uparrow \\ \text{ } \\ \downarrow \end{array} \begin{array}{l} -BD^{-1} \times \\ \text{ } \\ \times (-D^{-1}C) \end{array} \quad (\text{行基本変形})$$
$$\rightarrow \begin{pmatrix} A - BD^{-1}C & 0 \\ C & D \end{pmatrix}$$
$$\rightarrow \begin{pmatrix} M & 0 \\ 0 & D \end{pmatrix}.$$

これを基本変形行列を用いて書けば,

$$\begin{pmatrix} I_n & -BD^{-1} \\ 0 & I_m \end{pmatrix} \begin{pmatrix} A & B \\ C & D \end{pmatrix} \begin{pmatrix} I_n & 0 \\ -D^{-1}C & I_m \end{pmatrix} = \begin{pmatrix} M & 0 \\ 0 & D \end{pmatrix}.$$

$$\Leftrightarrow \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \begin{pmatrix} I_n & BD^{-1} \\ 0 & I_m \end{pmatrix} \begin{pmatrix} M & 0 \\ 0 & D \end{pmatrix} \begin{pmatrix} I_n & 0 \\ D^{-1}C & I_m \end{pmatrix}$$

両辺逆行列をとる.  $\downarrow (XY)^{-1} = Y^{-1}X^{-1}$

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} I_n & 0 \\ D^{-1}C & I_m \end{pmatrix}^{-1} \begin{pmatrix} M & 0 \\ 0 & D \end{pmatrix}^{-1} \begin{pmatrix} I_n & BD^{-1} \\ 0 & I_m \end{pmatrix}^{-1}$$

$$= \begin{pmatrix} I_n & 0 \\ -D^{-1}C & I_m \end{pmatrix} \begin{pmatrix} M^{-1} & 0 \\ 0 & D^{-1} \end{pmatrix} \begin{pmatrix} I_n & -BD^{-1} \\ 0 & I_m \end{pmatrix}$$

$$= \begin{pmatrix} M^{-1} & 0 \\ -D^{-1}CM^{-1} & D^{-1} \end{pmatrix} \begin{pmatrix} I_n & -BD^{-1} \\ 0 & I_m \end{pmatrix}$$

$$= \begin{pmatrix} M^{-1} & -M^{-1}BD^{-1} \\ -D^{-1}CM^{-1} & D^{-1} + D^{-1}CM^{-1}BD^{-1} \end{pmatrix}. \quad \blacksquare$$

• Th. 4.9.

$$\Lambda_{1,1} = (\Sigma_{1,1} - \Sigma_{1,2} \Sigma_{2,2}^{-1} \Sigma_{2,1})^{-1},$$

$$\Lambda_{1,2} = -\Lambda_{1,1} \Sigma_{1,2} \Sigma_{2,2}^{-1},$$

$$\Lambda_{2,2} = \Sigma_{2,2}^{-1} + \Sigma_{2,2}^{-1} \Sigma_{1,2} \Lambda_{1,1} \Sigma_{1,2} \Sigma_{2,2}^{-1}.$$

$\therefore \Lambda_{1,1}$  は,  $\mathbb{R}$  の Sherman-Morrison-Woodbury formula を求める.

Th. (Sherman-Morrison-Woodbury)

$A \in M_n(\mathbb{R})$ : regular,  $D \in M_m(\mathbb{R})$ : regular.

$B \in M_{n,m}(\mathbb{R})$ ,  $C \in M_{m,n}(\mathbb{R})$ ,  $D^{-1} + CA^{-1}B$ : regular

$$\Rightarrow (A + BDC)^{-1} = A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1}.$$

pf. ブロック行列の基本変形を考慮.

行の交換 ↓

← 列の交換

$$\begin{pmatrix} 0 & I_m \\ I_n & 0 \end{pmatrix} \begin{pmatrix} A & -B \\ C & D^{-1} \end{pmatrix} \begin{pmatrix} 0 & I_n \\ I_m & 0 \end{pmatrix} = \begin{pmatrix} D^{-1} & C \\ -B & A \end{pmatrix}.$$

両辺逆行列をとると,

$$\begin{pmatrix} 0 & I_n \\ I_m & 0 \end{pmatrix}^{-1} \begin{pmatrix} A & -B \\ C & D^{-1} \end{pmatrix}^{-1} \begin{pmatrix} 0 & I_m \\ I_n & 0 \end{pmatrix}^{-1} = \begin{pmatrix} D^{-1} & C \\ -B & A \end{pmatrix}^{-1}$$

$$\Leftrightarrow \begin{pmatrix} A & -B \\ C & D^{-1} \end{pmatrix}^{-1} = \begin{pmatrix} 0 & I_n \\ I_m & 0 \end{pmatrix} \begin{pmatrix} D^{-1} & C \\ -B & A \end{pmatrix}^{-1} \begin{pmatrix} 0 & I_m \\ I_n & 0 \end{pmatrix} \quad -(*)$$

上のTh.より,  $M = A + BDC$

$$\begin{pmatrix} A & -B \\ C & D^{-1} \end{pmatrix}^{-1} = \begin{pmatrix} M^{-1} & * \\ * & * \end{pmatrix}.$$

$$\begin{pmatrix} D^{-1} & C \\ -B & A \end{pmatrix}^{-1} = \begin{pmatrix} * & * \\ * & A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1} \end{pmatrix}$$

(\*) の左上ブロックを比較して,

$$(A + BDC)^{-1} = A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1}. \quad \square$$

Cor.  $P \in M_n(\mathbb{R}), R \in M_m(\mathbb{R}), P, R \succ 0 \quad B \in M_{m,n}(\mathbb{R})$

$$\Rightarrow (P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}.$$

pf.  $P, R \succ 0$  より  $R^{-1} + B P^{-1} B^T \succ 0$ .

SMW公理より,

$$(P^{-1} + B^T R^{-1} B)^{-1} = P - P B^T (R + B P B^T)^{-1} B P$$

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T R^{-1} - P B^T (B P B^T + R)^{-1} B P B^T R^{-1}.$$

$$(\text{r.h.s.}) = P B^T (B P B^T + R)^{-1} ((\cancel{B P B^T} + R) R^{-1} - \cancel{B P B^T} R^{-1})$$

$$= P B^T (B P B^T + R)^{-1} I_m$$

$$= P B^T (B P B^T + R).$$



- Gauss分布のpdfの対数をとると,

$$\log \mathcal{N}(x | \mu, \Sigma)$$

$$= -\frac{1}{2} (M \log 2\pi + \log \det \Sigma) - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)$$

$$= -\frac{1}{2} \left( (x - \mu)^T \Sigma^{-1} (x - \mu) + \log \det \Sigma + M \log 2\pi \right)$$

$$= -\frac{1}{2} \left( x^T \Sigma^{-1} x - 2x^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu + \log \det \Sigma + M \log 2\pi \right).$$

$$= -\frac{1}{2} \left( x^T \Sigma^{-1} x - 2x^T \Sigma^{-1} \mu + C(\mu, \Sigma) \right)$$

$\because C(\mu, \Sigma)$  とおく

この部分を覚えて Gauss分布の117x-9 p.100を参照.

- $p(x_1 | x_2)$  を求める.

$$\log p(x_1 | x_2)$$

$x_2$ : given のときは定数

$$= \log p(x_1, x_2) - \log p(x_2)$$

$$= -\frac{1}{2} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}^T \begin{pmatrix} \Lambda_{1,1} & \Lambda_{1,2} \\ \Lambda_{2,1} & \Lambda_{2,2} \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} + \text{Const.}$$

$$= -\frac{1}{2} \begin{pmatrix} (x_1 - \mu_1)^T \Lambda_{1,1} + (x_2 - \mu_2)^T \Lambda_{2,1} \\ (x_1 - \mu_1)^T \Lambda_{1,2} + (x_2 - \mu_2)^T \Lambda_{2,2} \end{pmatrix}^T \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} + \text{Const.}$$

$x_2$  は関係ない  
項は Const. に含む.

$$= -\frac{1}{2} \left( (x_1 - \mu_1)^T \Lambda_{1,1} (x_1 - \mu_1) + 2(x_2 - \mu_2)^T \Lambda_{2,1} (x_1 - \mu_1) \right) + \text{Const.}$$

$$= -\frac{1}{2} \left( x_1^T \Lambda_{1,1} x_1 - 2x_1^T (\Lambda_{1,1} \mu_1 - \Lambda_{1,2} (x_2 - \mu_2)) \right) + \text{Const.}$$

$$= -\frac{1}{2} \left( \mathbf{x}_1^T \underbrace{\Lambda_{1,1}}_{=: \Lambda_{1|2}} \mathbf{x}_1 - 2 \mathbf{x}_1^T \Lambda_{1,1} \left( \underbrace{\mu_1 - \Lambda_{1,1}^{-1} \Lambda_{1,2}}_{=: \mu_{1|2}} (\mathbf{x}_2 - \mu_2) \right) \right) + \text{Const.}$$

$$\therefore \Lambda_{1|2} = \Lambda_{1,1}, \quad \mu_{1|2} = \mu_1 - \Lambda_{1,1}^{-1} \Lambda_{1,2} (\mathbf{x}_2 - \mu_2) \quad \text{etc}$$

$$p(\mathbf{x}_1 | \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 | \mu_{1|2}, \Lambda_{1|2}^{-1}) \quad \text{etc}$$

•  $p(\mathbf{x}_2)$  を求める。

$$\log p(\mathbf{x}_2)$$

$$= \log p(\mathbf{x}_1, \mathbf{x}_2) - \log p(\mathbf{x}_1 | \mathbf{x}_2)$$

$$= -\frac{1}{2} \left( \begin{pmatrix} (\mathbf{x}_1 - \mu_1)^T \Lambda_{1,1} + (\mathbf{x}_2 - \mu_2)^T \Lambda_{2,1} \\ (\mathbf{x}_1 - \mu_1)^T \Lambda_{1,2} + (\mathbf{x}_2 - \mu_2)^T \Lambda_{2,2} \end{pmatrix}^T \begin{pmatrix} \mathbf{x}_1 - \mu_1 \\ \mathbf{x}_2 - \mu_2 \end{pmatrix} \right)$$

$$+ \frac{1}{2} (\mathbf{x}_1 - \mu_{1|2})^T \Lambda_{1|2} (\mathbf{x}_1 - \mu_{1|2}) + \text{Const.}$$

↳  $\mathbf{x}_2$  に関する項は  
Const. に含む。

$$= -\frac{1}{2} \left( (\mathbf{x}_2 - \mu_2)^T \Lambda_{2,2} (\mathbf{x}_2 - \mu_2) + 2 (\mathbf{x}_2 - \mu_2)^T \Lambda_{2,1} (\mathbf{x}_1 - \mu_1) + 2 \mu_{1|2}^T \Lambda_{1|2} \mathbf{x}_2 - \mu_{1|2}^T \Lambda_{1|2} \mu_{1|2} \right) + \text{Const.}$$

$$= -\frac{1}{2} \left( \mathbf{x}_2^T \Lambda_{2,2} \mathbf{x}_2 - 2 \mathbf{x}_2^T \Lambda_{2,2} \mu_2 + 2 \mathbf{x}_2^T \Lambda_{2,1} (\mathbf{x}_1 - \mu_1) \right)$$

$$+ 2 \left( \mu_1 - \Lambda_{1,1}^{-1} \Lambda_{1,2} (\mathbf{x}_2 - \mu_2) \right)^T \Lambda_{1,1} \mathbf{x}_1$$

$$- \left( \mu_1 - \Lambda_{1,1}^{-1} \Lambda_{1,2} (\mathbf{x}_2 - \mu_2) \right)^T \Lambda_{1,1} \left( \mu_1 - \Lambda_{1,1}^{-1} \Lambda_{1,2} (\mathbf{x}_2 - \mu_2) \right) \right) + \text{Const.}$$

$$= -\frac{1}{2} \left( \mathbf{x}_2^T (\Lambda_{2,2} - \Lambda_{2,1} \Lambda_{1,1}^{-1} \Lambda_{1,2}) \mathbf{x}_2 \right)$$

$$- 2 \mathbf{x}_2^T \left( \Lambda_{2,2} \mu_2 - \Lambda_{2,1} (\mathbf{x}_1 - \mu_1) + \Lambda_{2,1} \mathbf{x}_1 - \Lambda_{2,1} \mu_1 - \Lambda_{2,1} \Lambda_{1,1}^{-1} \Lambda_{1,2} \mu_2 \right) \right) + \text{Const.}$$

$$= -\frac{1}{2} \left( \mathbf{x}_2^T (\Lambda_{2,2} - \Lambda_{2,1} \Lambda_{1,1}^{-1} \Lambda_{1,2}) \mathbf{x}_2 - 2 \mathbf{x}_2^T (\Lambda_{2,2} - \Lambda_{2,1} \Lambda_{1,1}^{-1} \Lambda_{1,2}) \mu_2 \right) + \text{Const.}$$

∴ "SMW公式" を使うと,

$$\begin{aligned}
 & (\Lambda_{2,2} - \Lambda_{2,1} \Lambda_{1,1}^{-1} \Lambda_{1,2})^{-1} \quad \downarrow \text{SMW公式} \\
 & = \Lambda_{2,2}^{-1} + \Lambda_{2,2}^{-1} \Lambda_{2,1} (\Lambda_{1,1} - \Lambda_{1,2} \Lambda_{2,2}^{-1} \Lambda_{2,1})^{-1} \Lambda_{1,2} \Lambda_{2,2}^{-1} \\
 & = \Sigma_{2,2}.
 \end{aligned}$$

$\left( \begin{array}{cc} \Sigma_{1,1} & \Sigma_{1,2} \\ \Sigma_{2,1} & \Sigma_{2,2} \end{array} \right)$   
 $\left( \begin{array}{cc} \Lambda_{1,1} & \Lambda_{1,2} \\ \Lambda_{2,1} & \Lambda_{2,2} \end{array} \right)^{-1}$   
 ↓ ブロック行列 (右下のブロック)

$$\therefore p(x_2) = \mathcal{N}(x_2 | \mu_2, \Sigma_{2,2}).$$

・ 対称性を考え、

$$p(x_2 | x_1) = \mathcal{N}(x_2 | \mu_{2|1}, \Sigma_{2|1}),$$

$$\mu_{2|1} = \mu_2 - \Lambda_{2,2}^{-1} \Lambda_{2,1} (x_1 - \mu_1),$$

$$\Sigma_{2|1} = \Lambda_{2,2}^{-1}.$$

$$p(x_1) = \mathcal{N}(x_1 | \mu_1, \Sigma_{1,1}).$$

### 3.2.3 指数型分布族

#### 3.2.3.1 定義.

Def. (指数型分布族)

指数型分布族 (exponential family) とは、次のように pdf / pmf を表せる

確率分布の族:

$$p(x | \eta) = h(x) \exp(\eta^T \boldsymbol{\tau}(x) - a(\eta))$$

$\eta$ : 自然パラメータ (natural parameter),

$h(x)$ : 基底尺度 (base measure)

$a(\eta)$ : 対数分配記数関数 (log partition func.) という.

- $a(\eta)$  は  $p(x|\eta)$  が確率分布になるように調整するための項.

$$a(\eta) = \log \int h(x) \exp(\eta^T t(x)) dx.$$

$x$  が不明でも,  $t(x)$  の値が分かれば  $\Theta$  の最尤推定をするのに十分になる.

- $t(x)$  は分布のパラメータ  $\Theta$  の **十分統計量** (sufficient statistic) になっている.

cf.  $t(x)$  が  $\Theta$  の十分統計量とは,  $t(x)$  を与えたときの  $x$  の条件付き分布が

$\Theta$  に依存しないことという. 同値な表現では,  $x$ : given のときの

$\Theta$  の尤度関数  $L(\Theta|x)$  が,

$$L(\Theta|x) = g_{\Theta}(t(x)) \cdot h(x)$$

と表せることという. ( $g_{\Theta}$  は  $t$  の関数で  $\Theta$  を含む)

対数尤度  $l(\Theta|x) = \log g_{\Theta}(t(x)) + \log h(x)$  なのに最大法では

$$\frac{\partial}{\partial \Theta} l(\Theta|x) = \frac{\partial}{\partial \Theta} g_{\Theta}(t(x)) = 0$$

となる  $\Theta$  を求めることになり,  $x$  が分かると  $t(x)$  が分かれば  $\Theta$  の

最尤推定ができる.

### 3.2.3.2 分布の例.

- Bernoulli 分布は指数型分布族に入る.

$$\text{pf. Bern}(x|\mu) = \mu^x (1-\mu)^{1-x} = \exp(\log \mu^x (1-\mu)^{1-x})$$

$$= \exp(x \log \mu + (1-x) \log (1-\mu)) = \exp\left(x \log \frac{\mu}{1-\mu} + \log (1-\mu)\right)$$

$$\therefore h(x) = 1, \eta = \log \frac{\mu}{1-\mu}, t(x) = x, a(\eta) = \log (1+e^{\eta}).$$

$$\text{① } e^{\eta} = \frac{\mu}{1-\mu} \Leftrightarrow \mu = \frac{e^{\eta}}{1+e^{\eta}}$$

• Poisson分布 pmf:  $\text{Poi}(x|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$ .

Poisson分布は指数型分布族:  $\lambda$ .

pf.  $\text{Poi}(x|\lambda) = \frac{1}{x!} \exp(x \log \lambda - \lambda)$ .

$h(x) = \frac{1}{x!}$ ,  $\eta = \log \lambda$ ,  $t(x) = x$ ,  $a(\eta) = \lambda = e^\eta$ .

• Gauss分布や多項分布も指数型分布族 (証明略)

### 3.2.3.3 対数分配密度関数と十分統計量の関係.

•  $\nabla_\eta a(\eta) = \mathbb{E}[t(x)]$ .

pf.  $\nabla_\eta a(\eta)$

$= \nabla_\eta \log \int h(x) \exp(\eta^T t(x)) dx$

$= \left( \int h(x) \exp(\eta^T t(x)) dx \right)^{-1} \nabla_\eta \left( \int h(x) \exp(\eta^T t(x)) dx \right)$

$= \exp(-a(\eta)) \int h(x) \nabla_\eta \exp(\eta^T t(x)) dx$  ↓ 微分と積分の順序交換

$= \exp(-a(\eta)) \int t(x) h(x) \exp(\eta^T t(x)) dx$

$= \int t(x) \frac{h(x) \exp(\eta^T t(x) - a(\eta))}{p(x|\eta)} dx$

$= \mathbb{E}[t(x)]$ . ▣

•  $\frac{\partial^2}{\partial \eta \partial \eta^T} a(\eta) = V[t(x)]$

pf.  $\frac{\partial^2}{\partial \eta \partial \eta^T} a(\eta)$

$= \frac{\partial}{\partial \eta^T} \int t(x) p(x|\eta) dx$

$= \int t(x) \frac{\partial}{\partial \eta^T} p(x|\eta) dx$  ↑ 1つだけ ↓ 微分と積分の順序交換

$$\begin{aligned}
&= \int \boldsymbol{t}(x) h(x) \exp(\boldsymbol{\eta}^T \boldsymbol{t}(x) - a(\boldsymbol{\eta})) \left( \boldsymbol{t}(x)^T - \frac{\partial}{\partial \boldsymbol{\eta}^T} a(\boldsymbol{\eta}) \right) dx \\
&= \int \boldsymbol{t}(x) \boldsymbol{t}(x)^T p(x|\boldsymbol{\eta}) dx - \left( \int \boldsymbol{t}(x) p(x|\boldsymbol{\eta}) dx \right) \mathbb{E}[\boldsymbol{t}(x)]^T \\
&= \mathbb{E}[\boldsymbol{t}(x) \boldsymbol{t}(x)^T] - \mathbb{E}[\boldsymbol{t}(x)] \mathbb{E}[\boldsymbol{t}(x)]^T \\
&= V[\boldsymbol{t}(x)]. \quad \square
\end{aligned}$$

### 3.2.4 分布の共役性

Def. (共役事前分布)

指数型分布族  $p(x|\boldsymbol{\eta})$  に対し, 次の形の  $\boldsymbol{\eta}$  の事前分布  $p_\lambda(\boldsymbol{\eta})$  を  $p(x|\boldsymbol{\eta})$  の **共役事前分布** (conjugate prior) といい:

$$p_\lambda(\boldsymbol{\eta}) = h_c(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \boldsymbol{\lambda}_1 - a(\boldsymbol{\eta}) \boldsymbol{\lambda}_2 - a_c(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2)).$$

- $a_c(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2)$  は  $p_\lambda(\boldsymbol{\eta})$  が確率密度関数になるように正規化するための項

$$a_c(\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2) = \log \int h_c(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \boldsymbol{\lambda}_1 - a(\boldsymbol{\eta}) \boldsymbol{\lambda}_2) d\boldsymbol{\eta}.$$

#### 3.2.4.1 事後分布の解析的計算

- 共役事前分布を使うと, 指数型分布族の尤度関数に対し 事後分布も事前分布と同じ形式になる.

- $\mathcal{X} = \{x_1, \dots, x_N\}$ : Observed. 事後分布  $p_\lambda(\boldsymbol{\eta}|\mathcal{X})$  は,

$$\begin{aligned}
&p_\lambda(\boldsymbol{\eta}|\mathcal{X}) \\
&\propto p_\lambda(\boldsymbol{\eta}) p(\mathcal{X}|\boldsymbol{\eta}) \\
&= p_\lambda(\boldsymbol{\eta}) \prod_{n=1}^N p(x_n|\boldsymbol{\eta})
\end{aligned}$$

Bayesの定理  $p_\lambda(\boldsymbol{\eta}|\mathcal{X}) = \frac{p_\lambda(\boldsymbol{\eta}) p(\mathcal{X}|\boldsymbol{\eta})}{p(\mathcal{X})}$

↑ 事後分布      ↓ 事前分布      ← 尤度

↳  $\mathcal{X}$ : given data 定数.      周辺尤度

$$= h_c(\eta) \exp(\eta^T \lambda_1 - a(\eta) \lambda_2 - a_c(\lambda_1, \lambda_2)) \prod_{n=1}^N h(x_n) \exp(\eta^T t(x_n) - a(\eta))$$

$$\propto h_c(\eta) \exp(\eta^T (\lambda_1 + \sum_{n=1}^N t(x_n)) - a(\eta) (\lambda_2 + N)).$$

∴ 事後分布は事前分布と同じ形式.  $\eta$  は

$$\hat{\lambda}_1 = \lambda_1 + \sum_{n=1}^N t(x_n), \quad \hat{\lambda}_2 = \lambda_2 + N$$

と更新される. ← 事後分布が解析的に求まる!

### 3.2.4.2 予測分布の解析的計算.

•  $\eta$  が観測された後の未観測データ  $x_*$  の予測分布も解析的に求まる.

$$p(x_* | \mathcal{X}) = \mathbb{E}_{p(\eta | \mathcal{X})} [p(x_* | \eta)]$$

$$= \int p(x_* | \eta) p(\eta | \mathcal{X}) d\eta$$

$$= \int h(x_*) \exp(\eta^T t(x_*) - a(\eta)) h_c(\eta) \exp(\eta^T \hat{\lambda}_1 - a(\eta) \hat{\lambda}_2 - a_c(\hat{\lambda}_1, \hat{\lambda}_2)) d\eta.$$

$$= h(x_*) \exp(-a_c(\hat{\lambda}_1, \hat{\lambda}_2)) \int \frac{h(x_*) \exp(\eta^T (\hat{\lambda}_1 + t(x_*)) - a(\eta) (\hat{\lambda}_2 + 1))}{\exp(a_c(\hat{\lambda}_1 + t(x_*), \hat{\lambda}_2 + 1))} d\eta$$

$$= h(x_*) \exp(a_c(\hat{\lambda}_1 + t(x_*), \hat{\lambda}_2 + 1) - a_c(\hat{\lambda}_1, \hat{\lambda}_2)).$$

これは, 一般には指数型分布族には当てはまらない.

### 3.2.4.3 Bernoulli分布のパラメータの推論

• Bernoulli分布の共役事前分布は  $\Gamma$ -分布:

$$\text{pdf: } \text{Beta}(\mu | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} \mu^{\alpha-1} (1-\mu)^{\beta-1}, \quad \alpha, \beta > 0$$

$$\text{Bern}(x | \mu) = \exp\left(\frac{x \log \frac{\mu}{1-\mu}}{\eta} - \frac{\log(1+e^\eta)}{a(\eta)}\right).$$

自然パラメータ  $\eta = \log \frac{\mu}{1-\mu}$  を用いて  $\Gamma$ -分布を表現する.

$$e^\eta = \frac{\mu}{1-\mu} \Leftrightarrow (1-\mu)e^\eta = \mu \Leftrightarrow \mu = \frac{e^\eta}{1+e^\eta}$$

pdf の変数変換の式より  $\eta$  についての pdf  $\text{Beta}_\eta(\eta | \lambda_1, \lambda_2)$  は

$$\begin{aligned} & \text{Beta}_\eta(\eta | \lambda_1, \lambda_2) \\ &= \text{Beta}(\mu | \alpha, \beta) \left| \frac{d\mu}{d\eta} \right| \\ &= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha-1} (1-\mu)^{\beta-1} \left| \frac{e^\eta(1+e^\eta) - e^\eta \cdot e^\eta}{(1+e^\eta)^2} \right| \\ &= \exp\left( (\alpha-1)\log \mu + (\beta-1)\log(1-\mu) + \log \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} + \eta - 2 \frac{\log(1+e^\eta)}{a(\eta)} \right) \\ &= \exp\left( (\alpha-1)(\eta - \frac{\log(1+e^\eta)}{a(\eta)}) + (\beta-1)(-\frac{\log(1+e^\eta)}{a(\eta)}) + \log \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} + \eta - 2a(\eta) \right) \\ &= \exp\left( \eta\alpha - a(\eta)(\alpha+\beta) + \log \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \right). \end{aligned}$$

$$\therefore h_c(\eta) = 1, \lambda_1 = \alpha, \lambda_2 = \alpha + \beta, a_c(\lambda_1, \lambda_2) = \log \frac{\Gamma(\lambda_1)\Gamma(\lambda_2 - \lambda_1)}{\Gamma(\lambda_2)}$$

•  $\mathcal{X} = \{x_1, \dots, x_N\}$  ( $x_n \in \{0, 1\}$ ) 観測値とすると

石置 p1 = 事後事前分布!

$$\hat{\lambda}_1 = \lambda_1 + \sum_{n=1}^N t(x_n) = \alpha + \sum_{n=1}^N x_n =: \hat{\alpha}$$

$$\hat{\lambda}_2 = \lambda_2 + N = \alpha + \beta + N =: \hat{\alpha} + \hat{\beta}$$

$$\text{と更新される. } \hat{\beta} = \alpha + \beta + N - \hat{\alpha} = \beta + N - \sum_{n=1}^N x_n$$

$$\therefore \text{事後分布: } p(\mu | \mathcal{X}) = \text{Beta}(\mu | \hat{\alpha}, \hat{\beta})$$

予測分布は,

$$\begin{aligned} p(x_* | \mathcal{X}) &= \exp\left( a_c(\hat{\lambda}_1 + t(x_*), \hat{\lambda}_2 + 1) - a_c(\hat{\lambda}_1, \hat{\lambda}_2) \right) \\ &= \exp\left( \log \frac{\Gamma(\hat{\lambda}_1 + x_*)\Gamma(\hat{\lambda}_2 + 1 - \hat{\lambda}_1 - x_*)}{\Gamma(\hat{\lambda}_2 + 1)} \frac{\Gamma(\hat{\lambda}_2)}{\Gamma(\hat{\lambda}_1)\Gamma(\hat{\lambda}_2 - \hat{\lambda}_1)} \right) \\ &= \frac{\Gamma(\hat{\alpha} + x_*)\Gamma(\hat{\beta} + 1 - x_*)}{\hat{\lambda}_2 \Gamma(\hat{\lambda}_2)} \frac{\Gamma(\hat{\lambda}_2)}{\Gamma(\hat{\alpha})\Gamma(\hat{\beta})} \end{aligned}$$

$$= \frac{1}{\hat{\alpha} + \hat{\beta}} \frac{\Gamma(\hat{\alpha} + \alpha_*) \Gamma(\hat{\beta} + 1 - \alpha_*)}{\Gamma(\hat{\alpha}) \Gamma(\hat{\beta})}$$

$\Gamma(\hat{\alpha} + 1) = \hat{\alpha} \Gamma(\hat{\alpha})$

$$P(1 | \mathcal{X}) = \frac{1}{\hat{\alpha} + \hat{\beta}} \frac{\Gamma(\hat{\alpha} + 1) \Gamma(\hat{\beta} + 1 - 1)}{\Gamma(\hat{\alpha}) \Gamma(\hat{\beta})} = \frac{\hat{\alpha}}{\hat{\alpha} + \hat{\beta}}$$

$$P(0 | \mathcal{X}) = \frac{1}{\hat{\alpha} + \hat{\beta}} \frac{\Gamma(\hat{\alpha}) \Gamma(\hat{\beta} + 1)}{\Gamma(\hat{\alpha}) \Gamma(\hat{\beta})} = \frac{\hat{\beta}}{\hat{\alpha} + \hat{\beta}}$$

$$\therefore P(\alpha_* | \mathcal{X}) = \text{Bern}(\alpha_* | \frac{\hat{\alpha}}{\hat{\alpha} + \hat{\beta}})$$

cf.  $\Lambda^n$ -分布も指数型分布族.

$$\text{Beta}(\mu | \alpha, \beta) = \exp\left((\alpha-1)\log \mu + (\beta-1)\log(1-\mu) + \log \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}\right)$$

$$= \exp\left(\underbrace{\begin{pmatrix} \alpha-1 \\ \beta-1 \end{pmatrix}^T}_{\eta} \begin{pmatrix} \log \mu \\ \log(1-\mu) \end{pmatrix}_{t(\mu)} - \underbrace{\log \frac{\Gamma(\alpha-1+1)\Gamma(\beta-1+1)}{\Gamma(\alpha-1+\beta-1+2)}}_{a(\eta)}\right)$$

### 3.2.4.4 Gauss分布の精度パラメータの推論

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

平均パラメータ  $\mu$  を固定したとき, 精度パラメータ  $\gamma := \sigma^{-2}$  の共役事前分布は

ガンマ分布にたまる.

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

$$= \underbrace{\frac{1}{\sqrt{2\pi}}}_{h(x)} \exp\left(\underbrace{\gamma}_{\eta} \underbrace{\left(-\frac{1}{2}(x-\mu)^2\right)}_{t(x)} - \underbrace{\left(-\frac{1}{2}\log \gamma\right)}_{a(\eta)}\right)$$

$C_G(a, b) := \frac{b^a}{\Gamma(a)}$  (ガンマ分布の正規化定数) とする.

$$\text{Gam}(\eta | a, b) = C_G(a, b) \gamma^{a-1} e^{-b\gamma}$$

$$= \exp\left((a-1)\log \gamma - b\gamma + \log C_G(a, b)\right)$$

$$= \exp\left(\eta(-b) - \underbrace{\left(-\frac{1}{2}\log \gamma\right)}_{a(\eta)}(2(a-1)) - \left(-\log C_G(a, b)\right)\right)$$

$$= \exp\left(\frac{\eta(-b)}{\lambda_1} - a(\eta) \frac{2(a-1)}{\lambda_2} - \frac{(-\log C_G(a,b))}{a_c(\lambda_1, \lambda_2)}\right)$$

$$\therefore h_c(\eta) = 1, \quad \lambda_1 = -b, \quad \lambda_2 = 2(a-1).$$

$$a_c(\lambda_1, \lambda_2) = -\log C_G\left(\frac{\lambda_2}{2} + 1, -\lambda_1\right).$$

データ  $\mathcal{X} = \{x_1, \dots, x_N\}$  観測後の事後分布のパラメータは,

$$\hat{\lambda}_1 := -b + \sum_{n=1}^N t(x_n) = -\left(b + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2\right) =: -\hat{b}.$$

$$\hat{\lambda}_2 := \lambda_2 + N =: 2(\hat{a} - 1). \quad \Leftrightarrow \hat{a} = a + \frac{N}{2}.$$

$$\therefore \text{事後分布は } p(\gamma | \mathcal{X}) = \text{Gam}(\gamma | \hat{a}, \hat{b}).$$

予測分布は,

$$\begin{aligned} p(x_* | \mathcal{X}) &= \frac{1}{\sqrt{2\pi}} \exp\left(a_c(\hat{\lambda}_2 + t(x_*), \hat{\lambda}_2 + 1) - a_c(\hat{\lambda}_1, \hat{\lambda}_2)\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\log C_G\left(\frac{\hat{\lambda}_2 + 1}{2} + 1, -(\hat{\lambda}_2 + t(x_*))\right) + \log C_G\left(\frac{\hat{\lambda}_2}{2} + 1, -\hat{\lambda}_1\right)\right) \\ &= \frac{1}{\sqrt{2\pi}} \frac{C_G\left(\frac{\hat{\lambda}_2}{2} + 1, -\hat{\lambda}_1\right)}{C_G\left(\frac{\hat{\lambda}_2 + 1}{2} + 1, -(\hat{\lambda}_2 + t(x_*))\right)} \\ &= \frac{1}{\sqrt{2\pi}} \frac{C_G(\hat{a}, \hat{b})}{C_G\left(\hat{a} + \frac{1}{2}, \hat{b} - t(x_*)\right)} \\ &= \frac{1}{\sqrt{2\pi}} \frac{\Gamma(\hat{a} + \frac{1}{2})}{\Gamma(\hat{a})} \frac{\hat{b}^{\hat{a}}}{(\hat{b} - t(x_*))^{\hat{a} + \frac{1}{2}}} \\ &= \frac{1}{\sqrt{2\pi}} \frac{\Gamma(\hat{a} + \frac{1}{2})}{\Gamma(\hat{a})} \frac{\hat{b}^{\hat{a}}}{\hat{b}^{\hat{a} + \frac{1}{2}} \left(1 - \frac{1}{\hat{b}} t(x_*)\right)^{\hat{a} + \frac{1}{2}}} \\ &= \frac{1}{\sqrt{2\pi \hat{b}}} \frac{\Gamma(\hat{a} + \frac{1}{2})}{\Gamma(\hat{a})} \left(1 + \frac{1}{2\hat{b}} (x_* - \mu)^2\right)^{-\frac{2\hat{a} + 1}{2}} \\ &= \sqrt{\frac{\hat{a}/\hat{b}}{\pi \cdot 2\hat{a}}} \frac{\Gamma\left(\frac{2\hat{a} + 1}{2}\right)}{\Gamma(2\hat{a})} \left(1 + \frac{\hat{a}/\hat{b}}{2\hat{a}} (x_* - \mu)^2\right)^{-\frac{2\hat{a} + 1}{2}} \end{aligned}$$

=  $\mathcal{S}t(x_* | \mu_s, \lambda_s, \nu_s)$  : Student の  $t/\hat{\sigma}$  分布

$$\mu_s := \mu, \lambda_s := \frac{\hat{\sigma}^2}{b}, \nu_s := 2\hat{a}.$$

Student の  $t/\hat{\sigma}$  分布の pdf :  $\mathcal{S}t(x_* | \mu_s, \lambda_s, \nu_s) = \frac{\sqrt{\lambda_s}}{\sqrt{\pi \nu_s}} \frac{\Gamma(\frac{\nu_s+1}{2})}{\Gamma(\frac{\nu_s}{2})} \left(1 + \frac{\lambda_s}{\nu_s} (x_* - \mu_s)^2\right)^{-\frac{\nu_s+1}{2}}.$

•  $E_{\mathcal{S}t(\mu_s, \lambda_s, \nu_s)}[x] = \mu_s \quad (\nu_s > 1)$

•  $V_{\mathcal{S}t(\mu_s, \lambda_s, \nu_s)}[x] = \frac{\nu_s}{\lambda_s(\nu_s-2)} \quad (\nu_s > 2).$

} 証明は IT, 2, 3 下巻

## §3.3 Bayes 線形回帰.

### 3.3.1 モデル

データ:  $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$  連続値

$\mathcal{X} = \{x_1, \dots, x_N\}$ ,  $\mathcal{Y} = \{y_1, \dots, y_N\}$ .

$\mathcal{X}$  から  $\mathcal{Y}$  を予測する.

Bayes 線形回帰モデルの同時分布を次で定義

$$\begin{aligned} p(\mathcal{Y}, w | \mathcal{X}) &= p(w) p(\mathcal{Y} | \mathcal{X}, w) \\ &= p(w) \prod_{n=1}^N p(y_n | x_n, w). \end{aligned}$$

$y_n$  は分散  $\sigma_y^2$  の Gauss 分布に従うとする.

$$p(y_n | x_n, w) = \mathcal{N}(y_n | w^T \phi(x_n), \sigma_y^2),$$

$\phi: \mathbb{R}^{H_1} \rightarrow \mathbb{R}^{H_2}$ : 特徴量関数.

$w$  を学習する. Gauss 事前分布を与えておく.

$$p(w) = \mathcal{N}(w | 0, \sigma_w^2 I_{H_2})$$

### 3.3.2 学習と予測.

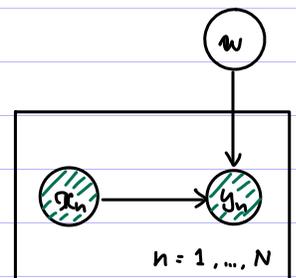
#### 3.3.2.1 事後分布の解析的計算.

事後分布は.

$$p(w | \mathcal{Y}, \mathcal{X}) = \frac{p(\mathcal{Y} | \mathcal{X}, w) p(w)}{p(\mathcal{Y} | \mathcal{X})} \propto p(\mathcal{Y} | \mathcal{X}, w) p(w)$$

←  $w$  に依存しない

対数をとって  $w$  について整理していく.





•  $\log p(w | y_*, x_*, y, \mathcal{X})$  について,

データ  $\{(x_1, y_1), \dots, (x_N, y_N), (x_*, y_*)\}$  : given とみたとし

事後分布を計算したものを考えたい。

$$\Sigma(x_*)^{-1} := \hat{\Sigma}^{-1} + \sigma_y^{-2} \phi(x_*) \phi(x_*)^T$$

$$\mu(y_*) := \Sigma(x_*) \left( \hat{\Sigma}^{-1} \hat{\mu} + \sigma_y^{-2} y_* \phi(x_*) \right)$$

とすれば,

$$\log p(w | y_*, x_*, y, \mathcal{X})$$

$$= \log \mathcal{N}(w | \mu(y_*), \Sigma(x_*))$$

$$= -\frac{1}{2} \left( H_1 \log 2\pi + \log \det \Sigma(x_*) + (w - \mu(y_*))^T \Sigma(x_*)^{-1} (w - \mu(y_*)) \right) + \text{const.}$$

$$= -\frac{1}{2} \left( \underbrace{w^T \Sigma(x_*)^{-1} w}_{y_* \text{ に依らず}} - 2 w^T \Sigma(x_*)^{-1} \mu(y_*) + \mu(y_*)^T \Sigma(x_*)^{-1} \mu(y_*) \right) + \text{const.}$$

$$= -\frac{1}{2} \left( \left( \hat{\mu}^T \hat{\Sigma}^{-1} + \sigma_y^{-2} y_* \phi(x_*)^T \right) \cancel{\Sigma(x_*)} \Sigma(x_*)^{-1} \Sigma(x_*) \left( \hat{\Sigma}^{-1} \hat{\mu} + \sigma_y^{-2} y_* \phi(x_*) \right) \right. \\ \left. - 2 w^T \cancel{\Sigma(x_*)}^{-1} \Sigma(x_*) \left( \hat{\Sigma}^{-1} \hat{\mu} + \sigma_y^{-2} y_* \phi(x_*) \right) \right) + \text{const.}$$

$$= -\frac{1}{2} \left( \underbrace{\hat{\mu}^T \hat{\Sigma}^{-1} \Sigma(x_*) \hat{\Sigma}^{-1} \hat{\mu}}_{y_* \text{ に依らず}} + 2 y_* \sigma_y^{-2} \phi(x_*)^T \Sigma(x_*) \hat{\Sigma}^{-1} \hat{\mu} \right. \\ \left. + y_*^2 \sigma_y^{-4} \phi(x_*)^T \Sigma(x_*) \phi(x_*) - 2 \left( \underbrace{w^T \hat{\Sigma}^{-1} \hat{\mu}}_{y_* \text{ に依らず}} + \sigma_y^{-2} y_* w^T \phi(x_*) \right) \right) + \text{const.}$$

$$= -\frac{1}{2} \left( \sigma_y^{-4} \phi(x_*)^T \Sigma(x_*) \phi(x_*) y_*^2 - 2 \sigma_y^{-2} \phi(x_*)^T (w - \Sigma(x_*) \hat{\Sigma}^{-1} \hat{\mu}) y_* \right) + \text{const.}$$

•  $\log p(y_* | x_*, y, \mathcal{X}, w)$

$\downarrow$   $y_*$  は  $y, \mathcal{X}$  に依らず。

$$= \log p(y_* | x_*, w)$$

$$= \log \mathcal{N}(y_* | w^T \phi(x_*), \sigma_y^2)$$

$$\begin{aligned}
&= -\frac{1}{2} \left( \log 2\pi + \log \sigma_y^2 + \sigma_y^{-2} (y_* - \mathbf{w}^T \phi(x_*))^2 \right) \\
&= -\frac{1}{2} \sigma_y^{-2} \left( y_*^2 - 2 \phi(x_*)^T \mathbf{w} y_* + \underbrace{(\mathbf{w}^T \phi(x_*))^2}_{y_* \text{ に依存しない}} \right) + \text{const.} \\
&= -\frac{1}{2} \sigma_y^{-2} (y_*^2 - 2 \phi(x_*)^T \mathbf{w} y_*) + \text{const.}
\end{aligned}$$

以上より,

$$\begin{aligned}
&\log p(y_* | x_*, \mathcal{Y}, \mathcal{X}) \\
&= -\frac{1}{2} \sigma_y^{-2} (y_*^2 - 2 \phi(x_*)^T \mathbf{w} y_*) \\
&\quad + \frac{1}{2} (\sigma_y^{-4} \phi(x_*)^T \Sigma(x_*) \phi(x_*) y_*^2 - 2 \sigma_y^{-2} \phi(x_*)^T (\mathbf{w} - \Sigma(x_*) \hat{\Sigma}^{-1} \hat{\mu}) y_*) + \text{const.} \\
&= -\frac{1}{2} \left( (\sigma_y^{-2} - \sigma_y^{-4} \phi(x_*)^T \Sigma(x_*) \phi(x_*)) y_*^2 - 2 \sigma_y^{-2} \phi(x_*)^T \Sigma(x_*) \hat{\Sigma}^{-1} \hat{\mu} y_* \right) + \text{const.}
\end{aligned}$$

とすると,

$$\begin{aligned}
&\log \mathcal{N}(x | \mu, \sigma^2) \\
&= -\frac{1}{2} \left( \log 2\pi + \log \sigma^2 + \sigma^{-2} (x - \mu)^2 \right) \quad \begin{array}{l} \downarrow x \text{ に関係ない } \epsilon \text{ は } 3.12 \\ \text{const. にまとめる} \end{array} \\
&= -\frac{1}{2} \left( \sigma^{-2} x^2 - 2\sigma^{-2} \mu x \right) + \text{const.}
\end{aligned}$$

ただし,

$$\begin{aligned}
(\sigma^2(x_*))^{-1} &:= \sigma_y^{-2} - \sigma_y^{-4} \phi(x_*)^T \Sigma(x_*) \phi(x_*), \\
\mu(x_*) &:= \sigma^2(x_*) \sigma_y^{-2} \phi(x_*)^T \Sigma(x_*) \hat{\Sigma}^{-1} \hat{\mu}
\end{aligned}$$

とすると,

$$p(y_* | x_*, \mathcal{Y}, \mathcal{X}) = \mathcal{N}(y_* | \mu(x_*), \sigma^2(x_*))$$

と証明

もう1回計算を繰り返す。

$$\begin{aligned}
 & \Sigma(x_*) \\
 &= \left( \hat{\Sigma}^{-1} + \sigma_y^{-2} \phi(x_*) \phi(x_*)^T \right)^{-1} \quad \text{Sherman-Morrison-Woodbury} \\
 &= \hat{\Sigma} - \hat{\Sigma} \phi(x_*) \left( (\sigma_y^{-2} I_1) + \phi(x_*)^T \hat{\Sigma} \phi(x_*) \right)^{-1} \phi(x_*)^T \hat{\Sigma} \\
 &= \hat{\Sigma} - \left( \sigma_y^2 + \phi(x_*)^T \hat{\Sigma} \phi(x_*) \right)^{-1} \hat{\Sigma} \phi(x_*) \phi(x_*)^T \hat{\Sigma} \quad \text{スカラー値}
 \end{aligned}$$

ゆえに,  $q = \phi(x_*)^T \hat{\Sigma} \phi(x_*)$  とおくと,

$$\begin{aligned}
 & (\sigma^2(x_*))^{-1} \\
 &= \sigma_y^{-2} - \sigma_y^{-4} \left( \phi(x_*)^T \hat{\Sigma} \phi(x_*) - (\sigma_y^2 + \phi(x_*)^T \hat{\Sigma} \phi(x_*))^{-1} (\phi(x_*)^T \hat{\Sigma} \phi(x_*))^2 \right) \\
 &= \sigma_y^{-2} - \sigma_y^{-4} \left( q - (\sigma_y^2 + q)^{-1} q^2 \right) \\
 &= \sigma_y^{-2} - \sigma_y^{-4} (\sigma_y^2 + q)^{-1} q (\cancel{\sigma_y^2 + q} \cancel{q}) \\
 &= \sigma_y^{-2} - \frac{\sigma_y^{-2} q}{\sigma_y^2 + q} \\
 &= \frac{1 + \cancel{\sigma_y^{-2} q} - \cancel{\sigma_y^{-2} q}}{\sigma_y^2 + q} \\
 &= (\sigma_y^2 + \phi(x_*)^T \hat{\Sigma} \phi(x_*))^{-1} \\
 \therefore \sigma^2(x_*) &= \sigma_y^2 + \phi(x_*)^T \hat{\Sigma} \phi(x_*).
 \end{aligned}$$

更に,

$$\begin{aligned}
 & \mu(x_*) \\
 &= \sigma^2(x_*) \sigma_y^{-2} \phi(x_*)^T \left( \hat{\Sigma} - (\sigma_y^2 + q)^{-1} \hat{\Sigma} \phi(x_*) \phi(x_*)^T \hat{\Sigma} \right) \hat{\Sigma}^{-1} \hat{\mu} \\
 &= \sigma^2(x_*) \sigma_y^{-2} \phi(x_*)^T \left( \hat{\mu} - \sigma^2(x_*)^{-1} \hat{\Sigma} \phi(x_*) \phi(x_*)^T \hat{\mu} \right)
 \end{aligned}$$

$$\begin{aligned}
&= \sigma^2(x_*) \sigma_y^{-2} \phi(x_*)^T \hat{\mu} - \sigma_y^{-2} q \phi(x_*)^T \hat{\mu} \\
&= \left( (\sigma_y^2 + q) \sigma_y^{-2} - \cancel{\sigma_y^{-2} q} \right) \phi(x_*)^T \hat{\mu} \\
&= \hat{\mu}^T \phi(x_*).
\end{aligned}$$

結局,

$$p(y_* | x_*, \mathcal{Y}, \mathcal{X}) = \mathcal{N}(y_* | \mu(x_*), \sigma^2(x_*)),$$

$$\mu(x_*) = \hat{\mu}^T \phi(x_*), \quad \sigma^2(x_*) = \sigma_y^2 + \phi(x_*)^T \hat{\Sigma} \phi(x_*).$$

### 3.3.2.3 最尤推定との比較.

- 最尤推定ではパラメータが変化しても大きく予測は変わらない

(外れ値などがなければ)

- Bayes推定では、パラメータが増えるにつれて予測の不確か性 (i.e., 予測値の分散) が変化していく.

### 3.3.3 周辺尤度.

- $\mathcal{X} = \{x_1, \dots, x_N\}$ ,  $\mathcal{Y} = \{y_1, \dots, y_N\}$  : given.

パラメータ  $w$  の同時分布  $p(\mathcal{Y}, w | \mathcal{X}) = p(\mathcal{Y} | \mathcal{X}, w) p(w)$  から積分除去した

$$p(\mathcal{Y} | \mathcal{X}) = \int p(\mathcal{Y} | \mathcal{X}, w) p(w) dw$$

を, Bayes線形回帰モデルの **周辺尤度** (marginal likelihood), ある意味

**エビデンス** (evidence) という.

モデルが与えられたもとでの、データ  $\mathcal{Y}$  の出現する尤もらしさを表す.

→ 複数個のモデルをエビデンスで定量的に比較できる。

・ 積分除去する代わりに、対数計算で  $p(y|\mathcal{X})$  を求める。

$$p(y, w | \mathcal{X}) = p(w | y, \mathcal{X}) p(y | \mathcal{X}) \quad \text{ただし,}$$

$$\begin{aligned} \log p(y | \mathcal{X}) &= \prod_{n=1}^N p(y_n | x_n, w) \\ &= \log p(w) + \log p(y | \mathcal{X}, w) - \log p(w | y, \mathcal{X}) \\ &= \log \mathcal{N}(w | 0, \sigma_w^2 I_{H_1}) + \sum_{n=1}^N \log \mathcal{N}(y_n | w^T \phi(x_n), \sigma_y^2) - \log \mathcal{N}(w | \hat{\mu}, \hat{\Sigma}) \\ &= -\frac{1}{2} \left( H_1 \log 2\pi + \log \det(\sigma_w^2 I_{H_1}) + w^T (\sigma_w^2 I_{H_1})^{-1} w \right) \\ &\quad - \frac{1}{2} \sum_{n=1}^N \left( \log 2\pi + \log \sigma_y^2 + \sigma_y^{-2} (y_n - w^T \phi(x_n))^2 \right) \\ &\quad + \frac{1}{2} \left( H_1 \log 2\pi + \log \det \hat{\Sigma} + (w - \hat{\mu})^T \hat{\Sigma}^{-1} (w - \hat{\mu}) \right) \\ &= -\frac{1}{2} \left( H_1 \log \sigma_w^2 + \sigma_w^{-2} w^T w \right. \\ &\quad \left. + N \log 2\pi + N \log \sigma_y^2 + \sigma_y^{-2} \sum_{n=1}^N y_n^2 - 2 w^T \left( \sigma_y^{-2} \sum_{n=1}^N y_n \phi(x_n) \right) + w^T \left( \sigma_y^{-2} \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \right) w \right) \\ &\quad - \left( \log \det \hat{\Sigma} - w^T \hat{\Sigma}^{-1} w + 2 w^T \hat{\Sigma}^{-1} \hat{\mu} - \hat{\mu}^T \hat{\Sigma}^{-1} \hat{\mu} \right) \\ &= -\frac{1}{2} \left( H_1 \log \sigma_w^2 + N \log 2\pi + N \log \sigma_y^2 + \sigma_y^{-2} \sum_{n=1}^N y_n^2 - \log \det \hat{\Sigma} - \hat{\mu}^T \hat{\Sigma}^{-1} \hat{\mu} \right. \\ &\quad \left. + w^T \left( \sigma_w^{-2} I_{H_1} + \sigma_y^{-2} \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \right) w - w^T \hat{\Sigma}^{-1} w \right) \\ &\quad = \hat{\Sigma}^{-1} \hat{\mu} \\ &= -\frac{1}{2} \left( \sigma_y^{-2} \sum_{n=1}^N y_n^2 + N \log \sigma_y^2 + N \log 2\pi + H_1 \log \sigma_w^2 - \log \det \hat{\Sigma} - \hat{\mu}^T \hat{\Sigma}^{-1} \hat{\mu} \right). \end{aligned}$$

$$\therefore p(y | \mathcal{X})$$

$$= \exp \left( -\frac{1}{2} \left( \sigma_y^{-2} \sum_{n=1}^N y_n^2 + N \log \sigma_y^2 + N \log 2\pi + H_1 \log \sigma_w^2 - \log \det \hat{\Sigma} - \hat{\mu}^T \hat{\Sigma}^{-1} \hat{\mu} \right) \right).$$

Bayes線形回帰ではこのように解析的に求める

一般にはエビデンスの計算は困難. 解析的に求めたとしても, 計算が  
大変なことも多い.

→ サンプリングにより積分を近似. 変分推論によるエビデンスの下界を  
求める. などして何とPIする.

### 3.3.4 逐次学習

逐次学習 (sequential learning) / オンライン学習 (online learning):

新規に入ってくる学習データに適応的に学習すること.

共役事前分布を使って解析的にパラメータ更新式が求まる場合は,

データの生成過程に順序の依存性を仮定しないなら,

逐次的にデータを与えるのと一度に全データを与える場合とで

最終的な事後分布はどちらも同じになる.

Bayes線形回帰の例. 最初のデータ  $\mathcal{D}_1 = \{x_1, y_1\}$ : given のときの事後分布.

$$p(w | y_1, x_1) \propto p(y_1 | w, x_1) p(w).$$

次のデータ  $\mathcal{D}_2 = \{x_2, y_2\}$ : given のとき,

$$p(w | y_1, x_1, y_2, x_2) \propto p(y_2 | w, x_2) p(w | y_1, x_1)$$

更新後の分布を  
事前分布として  
使う

モデルが複雑になると, 逐次学習の各更新で  $w$  の事後分布が解析的に求まらなくなる.

→ 近似的に事後分布を更新していく.

### 3.3.5 能動学習への応用.

- ・ 予測分布の計算

→ 予測対象の値に対する不確実性を、分散などの指標で評価できる.

→ 効率的に学習用ラベルデータを収集する **能動学習** (active learning)

に適用可能.

- ・ 能動学習:

ラベルの付いていないデータ  $\mathcal{X}_{\text{pool}}$  から、適切なデータ点  $\mathcal{X}_{\text{query}}$  を選ぶ.

そのデータにラベル  $y_{\text{query}}$  を要求する.

- ・  $\mathcal{X}_{\text{query}}$  の選び方 (一例):

新しい入力  $\mathcal{X}_*$  を与えたときの予測  $y_*$  の不確かなもの

(予測に自信がななもの) を選ぶ.

→ 予測分布の **エントロピー** (entropy) 最大の  $\mathcal{X}_* \in \mathcal{X}_{\text{pool}}$  を選ぶ:

$$\mathcal{X}_{\text{query}} \leftarrow \underset{\mathcal{X}_* \in \mathcal{X}_{\text{pool}}}{\operatorname{argmax}} F(\mathcal{X}_*),$$

$$F(\mathcal{X}_*) := -\mathbb{E}_{p(y_* | \mathcal{X}_*, y, \mathcal{X})} [\log p(y_* | \mathcal{X}_*, y, \mathcal{X})].$$

- ・ Bayes 線形回帰では,

$$F(\mathcal{X}_*)$$

$$= -\int p(y_* | \mathcal{X}_*, y, \mathcal{X}) \log p(y_* | \mathcal{X}_*, y, \mathcal{X}) dy_*$$

$$\begin{aligned}
&= - \int \mathcal{N}(y_* | \mu(\alpha_*), \sigma^2(\alpha_*)) \log \mathcal{N}(y_* | \mu(\alpha_*), \sigma^2(\alpha_*)) dy_* \\
&= \frac{1}{2} \int \mathcal{N}(y_* | \mu(\alpha_*), \sigma^2(\alpha_*)) \left( \log 2\pi + \log \sigma^2(\alpha_*) + \sigma^2(\alpha_*) (y_* - \mu(\alpha_*))^2 \right) dy_* \\
&= \frac{1}{2} \left( (\log 2\pi + \log \sigma^2(\alpha_*)) \underbrace{\int \mathcal{N}(y_* | \mu(\alpha_*), \sigma^2(\alpha_*)) dy_*}_{=1} \right. \\
&\quad \left. + \sigma^2(\alpha_*) \underbrace{\int (y_* - \mu(\alpha_*))^2 \mathcal{N}(y_* | \mu(\alpha_*), \sigma^2(\alpha_*)) dy_*}_{=\sigma^2(\alpha_*)} \right) \\
&= \frac{1}{2} (1 + \log \sigma^2(\alpha_*) + \log 2\pi).
\end{aligned}$$

とすると,  $\alpha_{\text{query}} \leftarrow \underset{\alpha_* \in \mathcal{X}_{\text{pool}}}{\text{argmax}} \sigma^2(\alpha_*)$  とする.

・ 同様にして未知函数  $f(\alpha)$  の最適値の探索ができる

(Bayes最適化)

・ 予測対象に対して弱い仮定を設定できる Gauss過程を用いる

・ 比較的一般的.

### 3.3.6 Gauss過程との関係.

・ 予測分布の分散は,

$$\begin{aligned}
&\sigma^2(\alpha_*) \\
&= \sigma_y^2 + \phi(\alpha_*)^T \hat{\Sigma} \phi(\alpha_*) \\
&= \sigma_y^2 + \phi(\alpha_*)^T \left( \sigma_y^{-2} \sum_{n=1}^N \phi(\alpha_n) \phi(\alpha_n)^T + \sigma_w^{-2} \mathbf{I}_{H_1} \right)^{-1} \phi(\alpha_*)
\end{aligned}$$

すなわち,

$$\Lambda = \sigma_w^{-2} \mathbf{I}_{H_1}, \quad \Phi = \underbrace{\begin{pmatrix} \phi(\alpha_1) & \dots & \phi(\alpha_N) \end{pmatrix}}_N \Bigg\} \begin{matrix} H_1 \\ \text{とすると,} \end{matrix}$$

$$\sigma^2(\alpha_*)$$

$$= \sigma_y^2 + \phi(\alpha_*)^T (\sigma_y^{-2} \Phi \Phi^T + \Lambda)^{-1} \phi(\alpha_*).$$

予測分布の平均は,  $\mathbf{y} = (y_1, \dots, y_N)^T$  とおくと,

$$\mu(\alpha_*)$$

$$= \hat{\beta}^T \phi(\alpha_*)$$

$$= \sigma_y^{-2} \left( \sum_{n=1}^N y_n \phi(\alpha_n)^T \right) (\sigma_y^{-2} \Phi \Phi^T + \Lambda)^{-1} \phi(\alpha_*)$$

$$= \sigma_y^{-2} \mathbf{y}^T \Phi^T (\sigma_y^{-2} \Phi \Phi^T + \Lambda)^{-1} \phi(\alpha_*)$$

$$= \sigma_y^{-2} \phi(\alpha_*)^T (\sigma_y^{-2} \Phi \Phi^T + \Lambda)^{-1} \Phi \mathbf{y}.$$

Sherman-Morrison-Woodbury の公式より,  $\Phi^T \Lambda^{-1} \Phi =: K$  とおくと,

$$(\sigma_y^{-2} \Phi \Phi^T + \Lambda)^{-1}$$

$$= \Lambda^{-1} - \Lambda^{-1} \Phi \left( (\sigma_y^{-2} \mathbf{I}_N)^{-1} + \underbrace{\Phi^T \Lambda^{-1} \Phi}_{=K} \right)^{-1} \Phi^T \Lambda^{-1}$$

$$= \Lambda^{-1} - \Lambda^{-1} \Phi (\sigma_y^2 \mathbf{I}_N + K)^{-1} \Phi^T \Lambda^{-1}$$

ゆえに,

$$\sigma^2(\alpha_*)$$

$$= \sigma_y^2 + \phi(\alpha_*)^T \left( \Lambda^{-1} - \Lambda^{-1} \Phi (\sigma_y^2 \mathbf{I}_N + K)^{-1} \Phi^T \Lambda^{-1} \right) \phi(\alpha_*)$$

$$= \sigma_y^2 + \phi(\alpha_*)^T \Lambda^{-1} \phi(\alpha_*) - \phi(\alpha_*)^T \Lambda^{-1} \Phi (\sigma_y^2 \mathbf{I}_N + K)^{-1} \Phi^T \Lambda^{-1} \phi(\alpha_*).$$

$$\mu(\alpha_*)$$

$$= \sigma_y^{-2} \phi(\alpha_*)^T \left( \Lambda^{-1} - \Lambda^{-1} \Phi (\sigma_y^2 \mathbf{I}_N + K)^{-1} \Phi^T \Lambda^{-1} \right) \Phi \mathbf{y} \quad \lambda^{-1} \epsilon'' \ll \epsilon.$$

$$= \sigma_y^{-2} \phi(\alpha_*)^T \Lambda^{-1} \left( I_{H_2} - \Phi (\sigma_y^2 I_N + K)^{-1} \Phi^T \Lambda^{-1} \right) \Phi y$$

↓  $\Phi \Sigma^{-1} \Phi^T$

$$= \sigma_y^{-2} \phi(\alpha_*)^T \Lambda^{-1} \left( \Phi - \Phi (\sigma_y^2 I_N + K)^{-1} \underbrace{\Phi^T \Lambda^{-1} \Phi}_{=K} \right) y$$

↓  $\Phi \Sigma^{-1} \Phi^T$

$$= \sigma_y^{-2} \phi(\alpha_*)^T \Lambda^{-1} \Phi \left( I_N - (\sigma_y^2 I_N + K)^{-1} K \right) y$$

$$= \sigma_y^{-2} \phi(\alpha_*)^T \Lambda^{-1} \Phi \left( I_N - (\sigma_y^2 I_N + K)^{-1} (\sigma_y^2 I_N + K - \sigma_y^2 I_N) \right) y$$

$$= \cancel{\sigma_y^{-2}} \phi(\alpha_*)^T \Lambda^{-1} \Phi \left( \cancel{I_N} - \cancel{I_N} + \cancel{\sigma_y^2} (\sigma_y^2 I_N + K)^{-1} \right) y$$

$$= \phi(\alpha_*)^T \Lambda^{-1} \Phi (\sigma_y^2 I_N + K)^{-1} y$$

特徴量函数  $\phi$  は, 二つの入力データ  $\alpha, \alpha'$  に対し常に:

$$k(\alpha, \alpha') := \phi(\alpha)^T \Lambda^{-1} \phi(\alpha') \quad (\text{カーネル函数}) \quad \leftarrow \begin{array}{l} \text{二つのデータ } \alpha, \alpha' \text{ の} \\ \text{“類似度” のようなもの} \end{array}$$

の形で現れているので,  $\Phi$  を設計しなくても  $K$  を直接設計すればよい (カーネルトリック):

$$K = \Phi^T \Lambda^{-1} \Phi$$

$$= \begin{pmatrix} \phi(\alpha_1)^T \Lambda^{-1} \\ \vdots \\ \phi(\alpha_N)^T \Lambda^{-1} \end{pmatrix} \begin{pmatrix} \phi(\alpha_1) & \dots & \phi(\alpha_N) \end{pmatrix}$$

$$= \begin{pmatrix} \phi(\alpha_1)^T \Lambda^{-1} \phi(\alpha_1) & \dots & \phi(\alpha_1)^T \Lambda^{-1} \phi(\alpha_N) \\ \vdots & \ddots & \vdots \\ \phi(\alpha_N)^T \Lambda^{-1} \phi(\alpha_1) & \dots & \phi(\alpha_N)^T \Lambda^{-1} \phi(\alpha_N) \end{pmatrix}$$

$$= (k(\alpha_i, \alpha_j))_{ij}$$

$$\phi(\alpha_*)^T \Lambda^{-1} \Phi$$

$$= \begin{pmatrix} \phi(\alpha_*)^T \Lambda^{-1} \phi(\alpha_1) \\ \vdots \\ \phi(\alpha_*)^T \Lambda^{-1} \phi(\alpha_N) \end{pmatrix} = \begin{pmatrix} k(\alpha_*, \alpha_1) \\ \vdots \\ k(\alpha_*, \alpha_N) \end{pmatrix}$$

## §3.4 最尤推定, MAP推定との関係

### 3.4.1 最尤推定と誤差最小化.

Th. 回帰モデルを,  $w \in \mathbb{R}^d$  かつ,  $\varepsilon \sim \mathcal{N}(0, \sigma_y^2)$  とし,

$$y = f(x; w) + \varepsilon$$

とする.  $N$ 個のデータ  $\mathcal{D} = \{x, y\}$  が与えられたときの

パラメータ  $w$  の最尤推定値  $w_{ML}$  は, 誤差関数

$$E(w) = \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n; w))^2$$

← 最尤推定は  
最小二乗法に対応する.

を最小化する  $w$  と一致する.

pf.  $\varepsilon = y - f(x; w) \sim \mathcal{N}(0, \sigma_y^2)$  より,

$$y_n \sim \mathcal{N}(f(x_n; w), \sigma_y^2)$$

である.  $\mathcal{D}$ : given のとき, モデルの尤度関数は,

$$p(y|x, w) = \prod_{n=1}^N p(y_n|x_n, w) = \prod_{n=1}^N \mathcal{N}(y_n | f(x_n; w), \sigma_y^2)$$

最尤推定値は,

$$\begin{aligned} w_{ML} & \quad \log \varepsilon \text{ と } \sigma \text{ の大小関係は変化する.} \\ & \quad \curvearrowright \\ &= \operatorname{argmax}_w p(y|x, w) = \operatorname{argmax}_w \log p(y|x, w) \\ &= \operatorname{argmax}_w \left[ \underbrace{-\frac{1}{2} \left( N \log 2\pi + N \log \sigma_y^2 + \frac{1}{\sigma_y^2} \sum_{n=1}^N (y_n - f(x_n; w))^2 \right)}_{\substack{\uparrow \text{定数は最大化に} \\ \text{関係ない} \rightarrow}} \right] \quad \downarrow \operatorname{argmin} \text{ に} \\ &= \operatorname{argmin}_w \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n; w))^2 \\ &= \operatorname{argmin}_w E(w). \end{aligned}$$



- $f(x; \omega)$  が NN のように複雑でパラメータの解析解が求まらない場合、勾配降下法により最適化する。

$$\begin{aligned} & \nabla_{\omega} \log p(y | \mathcal{X}, \omega) \\ &= -\sigma_y^{-2} \nabla_{\omega} \left( \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n; \omega))^2 \right) \\ &= -\sigma_y^{-2} \nabla_{\omega} E(\omega) \end{aligned}$$

なので、 $\log p(y | \mathcal{X}, \omega)$  の最大化の更新式は  $\alpha$  を学習率として

$$\begin{aligned} \omega &\leftarrow \omega + \alpha \nabla_{\omega} \log p(y | \mathcal{X}, \omega) \\ &= \omega - \alpha \sigma_y^{-2} \nabla_{\omega} E(\omega) \end{aligned}$$

$\alpha := \alpha' \sigma_y^{-2}$  とおくと、 $\omega \leftarrow \omega - \alpha \nabla_{\omega} E(\omega)$  となる、

$E$  を最小化する勾配降下法の更新式と同じ式が得られる。

### 3.4.2 MAP推定と正則化

- **最大事後確率推定** (maximum a posteriori estimation: MAP推定)

モデル  $y = f(x; \omega)$  のパラメータ  $\omega$  と、データ  $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$ : given のときの  $\omega$  の事後分布を最大化する  $\omega$  (モード) を推定する:

$$\omega_{\text{MAP}} := \underset{\omega}{\text{argmax}} p(\omega | \mathcal{Y}, \mathcal{X}) = \underset{\omega}{\text{argmax}} \log p(\omega | \mathcal{Y}, \mathcal{X}).$$

Th. 回帰モデル  $y = f(x; \omega) + \varepsilon$ ,  $\omega$  はパラメータ,  $\varepsilon \sim \mathcal{N}(0, \sigma_y^2)$  とし、

$$y = f(x; \omega) + \varepsilon$$

とする。また、 $\omega$  の事前分布を  $p(\omega) = \mathcal{N}(\omega | 0, \sigma_w^2 I_{H_2})$  とする。

$N$ 個のデータ  $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$  が与えられたときの

パラメータ  $w$  の MAP 推定値  $w_{\text{MAP}}$  は、コスト関数

$$J(w) = E(w) + \lambda \Omega_L(w) = \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n; w))^2 + \frac{\lambda}{2} \|w\|^2$$

を最小化する  $w$  と一致する。

pf.  $\log p(w | \mathcal{Y}, \mathcal{X})$  Bayes の定理

$$= \log p(\mathcal{Y} | \mathcal{X}, w) + \log p(w) + \text{const.}$$

$$= -\frac{1}{2} \sigma_y^{-2} \sum_{n=1}^N (y_n - f(x_n; w))^2$$

$$- \frac{1}{2} \left( H_1 \log 2\pi + \log \det(\sigma_w^2 I_{H_1}) + w^T (\sigma_w^2 I_{H_1})^{-1} w \right) + \text{const.}$$

$$= -\sigma_y^{-2} \left( \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n; w))^2 + \frac{1}{2} \frac{\sigma_w^{-2}}{\sigma_y^{-2}} \|w\|^2 \right) + \text{const.}$$

$$= -\sigma_y^{-2} \left( E(w) + \lambda \Omega_L(w) \right) + \text{const.}$$

$$\therefore w_{\text{MAP}} = \underset{w}{\text{argmax}} \log p(w | \mathcal{Y}, \mathcal{X})$$

$$= \underset{w}{\text{argmin}} J(w). \quad \square$$

事前分布  $\Sigma$  Laplace 分布  $\text{Lap}(w | \mu, b) = \frac{1}{2b} \exp\left(-\frac{1}{b} |w - \mu|\right)$

$$\text{よって } p(w) = \prod_{h=1}^{H_1} \text{Lap}(w_h | 0, \frac{1}{\lambda})$$

とすると、正則化項は  $\lambda \Omega_{L_1}(w) = \lambda \sum_{h=1}^{H_1} |w_h|$  となる:

$$\log p(w)$$

$$= H_1 \log \frac{\lambda}{2} - \sum_{h=1}^{H_1} \lambda |w_h|$$

$$= -\lambda \Omega_{L_1}(w) + \text{const.}$$

・ 事前分布に Gauss 分布や Laplace 分布を導入した MAP 推定は

最尤推定より過剰適合にロバスト (正則化の効果)

→ 利用するのは  $w_{MAP}$  のみ (点推定: point estimation) なので,

予測の不確実性の定量化やエビデンスによるモデルの評価

などはできない.

### 3.4.3 分類モデルに対する誤差関数.

#### 3.4.3.1 2値分類の場合.

・ ラベル:  $y_n \in \{0, 1\}$  とする.

→  $y_n \sim \text{Bern}(\mu_n)$  と仮定.

$\mu_n$  についても, 回帰モデル  $\eta_n = f(x_n; w) + \varepsilon_n$ ,  $\varepsilon_n \sim \mathcal{N}(0, \sigma_\eta^2)$

を利用して  $\mu_n = \text{Sig}(\eta_n)$  と仮定.

つまり,

$$\mu_n = \text{Sig}(\eta_n), \quad \eta_n \sim \mathcal{N}(f(x_n; w), \sigma_\eta^2)$$

・ 特:  $f(x_n; w) = w^T \phi(x_n)$  のとき ロジスティック回帰モデル.

・  $\eta = \text{Sig}^{-1}(\mu)$ .  $\text{Sig}(\eta) = \frac{1}{1+e^{-\eta}} = \mu$  なのぞ,

$\eta = \log \frac{\mu}{1-\mu}$ . ← Bernoulli 分布の自然パラメータ表現

シグモイド関数で Bernoulli 分布のパラメータを自然パラメータに変換.

・ データ  $\mathcal{D} = \{x, y\}$ : given とする. 対数尤度関数は.

$$\begin{aligned}
& \log p(y | \mathcal{X}, \omega) \\
&= \log \prod_{n=1}^N p(y_n | x_n, \omega) \\
&= \sum_{n=1}^N \log (\mu_n^{y_n} (1-\mu_n)^{1-y_n}) \\
&= \sum_{n=1}^N (y_n \log \mu_n + (1-y_n) \log (1-\mu_n))
\end{aligned}$$

→ このモデルに対する最尤法は、クロスエントロピー-誤差

$$-\sum_{n=1}^N (y_n \log \mu_n + (1-y_n) \log (1-\mu_n))$$

の最小化になる。

・ とは32". NNのとは32"は,  $\mu_n = f(x_n; \omega)$  と

決定的に求めたい。このときのクロスエントロピー-誤差は、 $\omega$ の関数として

$$\begin{aligned}
& E(\omega) \\
&= -\sum_{n=1}^N (y_n \log \text{Sig}(f(x_n; \omega)) + (1-y_n) \log (1 - \text{Sig}(f(x_n; \omega))))
\end{aligned}$$

と書ける。一方で、このモデルでは  $\mu_n$  は確率変数なので、

$\mu_n$  は  $\omega$ の関数にたよるから  $\omega$ で微分できない。最尤法を使うためには、

誤差が  $\omega$ の関数にたよるというわけではないのではないか？

→ 実は、上のクロスエントロピー-誤差も  $\omega$ の関数  $\tilde{E}(\omega)$  とみなせる。

$$\mu_n = f(x_n; \omega) + \varepsilon_n, \quad \varepsilon_n \sim \mathcal{N}(0, \sigma_\eta^2) \text{ なの2"},$$

$$\tilde{E}(\omega)$$

$$= -\sum_{n=1}^N (y_n \log \text{Sig}(f(x_n; \omega) + \varepsilon_n) + (1-y_n) \log (1 - \text{Sig}(f(x_n; \omega) + \varepsilon_n)))$$

これは  $w$  について微分がとれる。(reparameterization trick)

とはいえ、 $\epsilon_n$  が確率変数であることに変わりはないので、 $\tilde{E}(w)$  も

確率変数になる。そこで、平均的に loss が小さくなるように

$\tilde{E}(w)$  の期待値  $\mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\tilde{E}(w)]$  を最小化するようにする。

$$\begin{aligned} & \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\tilde{E}(w)] \\ &= - \sum_{n=1}^N \left( y_n \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\log \text{sig}(f(x_n; w) + \epsilon_n)] \right. \\ & \quad \left. + (1 - y_n) \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\log(1 - \text{sig}(f(x_n; w) + \epsilon_n))] \right) \end{aligned}$$

この最小化のためには  $\mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\tilde{E}(w)]$  の  $w$  に関する勾配を考えればよい

$$\begin{aligned} & \nabla_w \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\log \text{sig}(f(x_n; w) + \epsilon_n)] \quad \downarrow \text{微分と積分の順序交換.} \\ &= \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\nabla_w \log \text{sig}(f(x_n; w) + \epsilon_n)] \quad \downarrow \text{合成関数の微分. } \text{sig}'(z) = \text{sig}(z)(1 - \text{sig}(z)). \\ &= \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)} \left[ \frac{1}{\cancel{\text{sig}(f(x_n; w) + \epsilon_n)}} \cancel{\text{sig}(f(x_n; w) + \epsilon_n)} (1 - \text{sig}(f(x_n; w) + \epsilon_n)) \nabla_w f(x_n; w) \right] \\ &= \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)} [1 - \text{sig}(f(x_n; w) + \epsilon_n)] \nabla_w f(x_n; w), \\ & \nabla_w \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\log(1 - \text{sig}(f(x_n; w) + \epsilon_n))] \quad \downarrow \text{微分と積分の順序交換.} \\ &= \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)} \left[ \frac{-1}{\cancel{1 - \text{sig}(f(x_n; w) + \epsilon_n)}} \cancel{\text{sig}(f(x_n; w) + \epsilon_n)} (1 - \cancel{\text{sig}(f(x_n; w) + \epsilon_n)}) \nabla_w f(x_n; w) \right] \\ &= - \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)} [\text{sig}(f(x_n; w) + \epsilon_n)] \nabla_w f(x_n; w) \end{aligned}$$

よって、

$$\nabla_w \mathbb{E}_{\mathcal{N}(0, \sigma_z^2)}[\tilde{E}(w)]$$

$$\begin{aligned}
&= - \sum_{n=1}^N \left( y_n \mathbb{E}_{\mathcal{N}(0, \sigma_\varepsilon^2)} \left[ 1 - \text{Sig}(f(x_n; \omega) + \varepsilon_n) \right] \nabla_\omega f(x_n; \omega) \right. \\
&\quad \left. - (1 - y_n) \mathbb{E}_{\mathcal{N}(0, \sigma_\varepsilon^2)} \left[ \text{Sig}(f(x_n; \omega) + \varepsilon_n) \right] \nabla_\omega f(x_n; \omega) \right) \\
&= - \sum_{n=1}^N \left( y_n - \mathbb{E}_{\mathcal{N}(0, \sigma_\varepsilon^2)} \left[ \text{Sig}(f(x_n; \omega) + \varepsilon_n) \right] \right) \nabla_\omega f(x_n; \omega).
\end{aligned}$$

∴ z',

$$\begin{aligned}
&\mathbb{E}_{\mathcal{N}(0, \sigma_\varepsilon^2)} \left[ \text{Sig}(f(x_n; \omega) + \varepsilon_n) \right] \quad \downarrow \text{Taylor 展開. } \text{L2も良いかは不明} \dots \\
&= \mathbb{E}_{\mathcal{N}(0, \sigma_\varepsilon^2)} \left[ \sum_{i=0}^{\infty} \frac{1}{i!} \text{Sig}^{(i)}(f(x_n; \omega)) \varepsilon_n^i \right] \\
&= \sum_{i=0}^{\infty} \frac{1}{i!} \text{Sig}^{(i)}(f(x_n; \omega)) \underbrace{\mathbb{E}_{\mathcal{N}(0, \sigma_\varepsilon^2)} [\varepsilon_n^i]}_{\text{正規分布のモーメント. } i: \text{ odd } \rightarrow 0, i: \text{ even } \rightarrow (i-1)!! \sigma_\varepsilon^i} \\
&= \text{Sig}(f(x_n; \omega)) + \sum_{i=1}^{\infty} \frac{1}{(2i)!} \text{Sig}^{(2i)}(f(x_n; \omega)) (2i-1)!! \sigma_\varepsilon^{2i} \\
&\approx \text{Sig}(f(x_n; \omega))
\end{aligned}$$

∴ z'',  $\nabla_\omega \mathbb{E}_{\mathcal{N}(0, \sigma_\varepsilon^2)} [\tilde{E}(\omega)] \approx - \sum_{n=1}^N \left( y_n - \text{Sig}(f(x_n; \omega)) \right) \nabla_\omega f(x_n; \omega).$

一方 z''.

$$\begin{aligned}
&\nabla_\omega E(\omega) \\
&= - \sum_{n=1}^N \left( y_n \frac{\cancel{\text{Sig}(f(x_n; \omega))} (1 - \text{Sig}(f(x_n; \omega)))}{\text{Sig}(f(x_n; \omega))} \nabla_\omega f(x_n; \omega) \right. \\
&\quad \left. - (1 - y_n) \frac{\text{Sig}(f(x_n; \omega)) (1 - \cancel{\text{Sig}(f(x_n; \omega))})}{1 - \text{Sig}(f(x_n; \omega))} \nabla_\omega f(x_n; \omega) \right) \\
&= - \sum_{n=1}^N \left( y_n - \text{Sig}(f(x_n; \omega)) \right) \nabla_\omega f(x_n; \omega)
\end{aligned}$$

となるので, このモデルに対する最尤推定のパラメータは,

近似的には決定的に求める場合と同じになる.

### 3.4.3.2 多値分類の場合.

多値分類でも同様. D値分類と可.

$$y_n \in \{0,1\}^D, \sum_{d=1}^D y_{n,d} = 1 \text{ と可 (one-hot)}$$

$y_n \sim \text{Cat}(\pi_n)$  と仮定. さらに,

$$\pi_n = \text{Softmax}(\eta_n),$$

$$\eta_n = f(x_n; W) + \epsilon_n, \epsilon_n \sim \mathcal{N}(0, \sigma_\eta^2 I_D)$$

と仮定. 対数尤度関数は,

$$\begin{aligned} & \log p(y | x, W) \\ &= \log \prod_{n=1}^N p(y_n | x_n, W) \\ &= \sum_{n=1}^N \log p(y_n | x_n, W) \\ &= \sum_{n=1}^N \log \prod_{d=1}^D \pi_{n,d}^{y_{n,d}} \\ &= \sum_{n=1}^N \sum_{d=1}^D y_{n,d} \log \pi_{n,d} \end{aligned}$$

→ 最尤法は, クロスエントロピー誤差

$$- \sum_{n=1}^N \sum_{d=1}^D y_{n,d} \log \pi_{n,d}$$

の最小化になる.

・3.4.3.1 と同様.  $\eta_n$  を決定的に求める場合と, このモデルでのように確率的に決める場合は, 近似的には同じ (はず).