

# PRML輪読会

## 第4章線形識別モデル

# 目次

## 4.3 確率的識別モデル

### 4.3.1 固定基底関数

### 4.3.2 ロジスティック回帰

### 4.3.3 反復再重み付け最小二乗

### 4.3.4 多クラスロジスティック回帰

### 4.3.5 プロビット回帰

### 4.3.6 正準連結関数

## 4.4 ラプラス近似

### 4.4.1 モデルの比較とBIC

## 4.5 ベイズロジスティック回帰

### 4.5.1 ラプラス近似

### 4.5.2 予測分布

# 決定問題を解く3つのアプローチ

## 識別関数

入力データをいくつかのクラスの1つに割り当てる関数(識別関数)を直接導出する(確率は登場しない)  
(4.1で議論済)

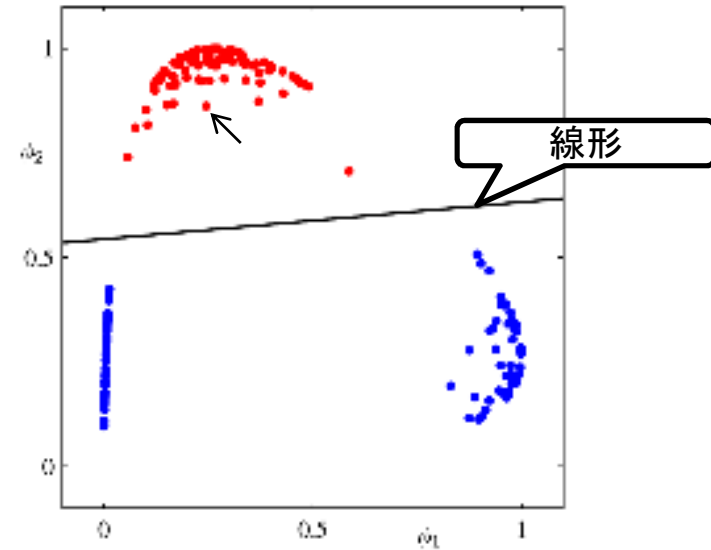
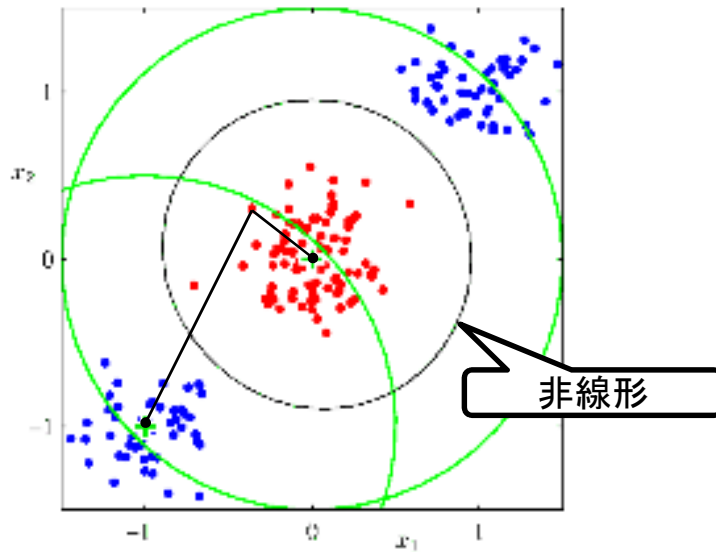
## 識別モデル

条件付き確率モデル $p(C_k | x)$ を、データから直接モデル化する  
(4.3で議論)

## 生成モデル

条件付き確率モデル $p(C_k | x)$ を、 $p(x|C_k)$ と $p(C_k)$ からベイズの定理を使って導出する(4.2で議論済)

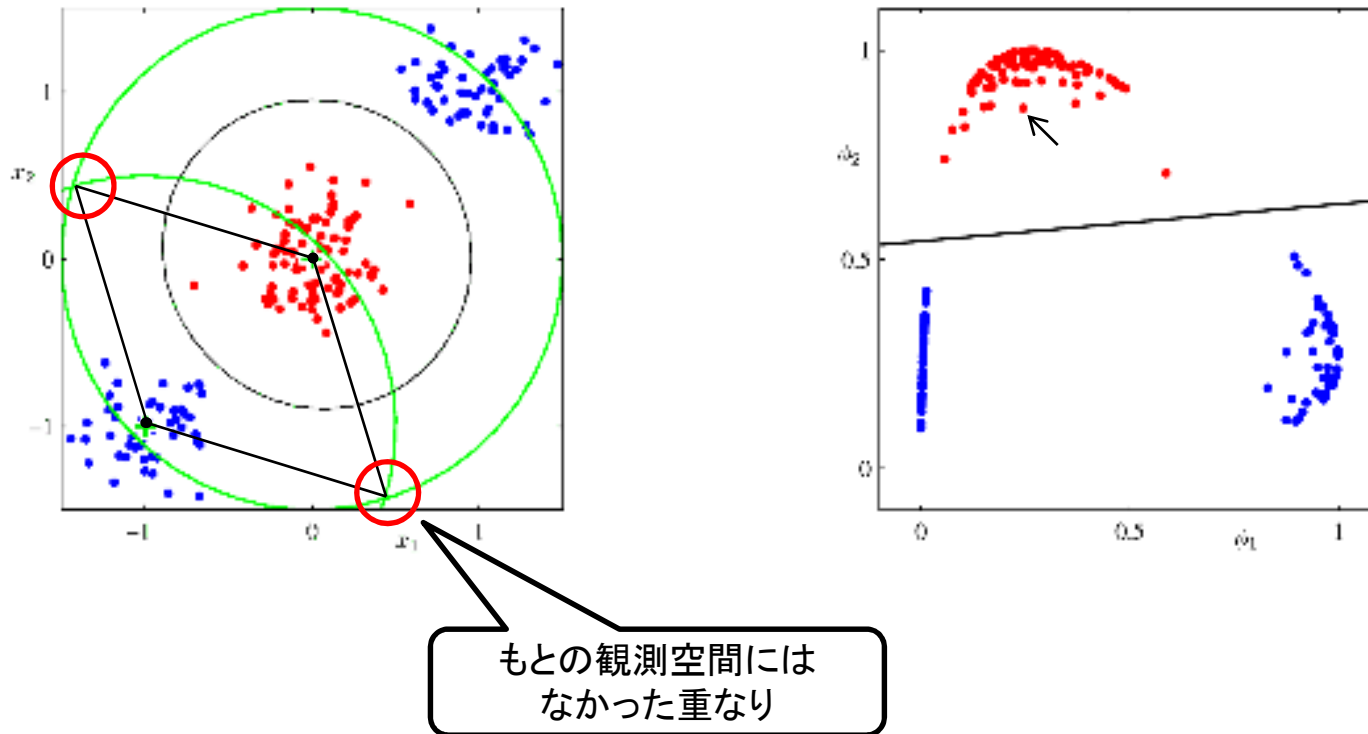
### 4.3.1 固定基底関数



$$\begin{aligned}\phi_1(x) &= \exp\left(-\left\|x - \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right\|\right) \\ \phi_2(x) &= \exp(-\|x\|)\end{aligned}\quad (\text{たぶんこんな感じ})$$

あらかじめ非線形変換を行っておけば、変換後の特徴空間においては線形決定境界を得られるかもしれない

### 4.3.1 固定基底関数



非線形変換を行うことで、もとの空間にはなかった重なりを作ってしまうこともある。

とはいえ有用なので固定基底関数使います。

## 4.3.2 ロジスティック回帰

ロジスティック回帰とは

(2クラス分類問題において) 以下の式であらわされるモデル

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (4.87)$$

$$\text{ただし、}\sigma(a) = \frac{1}{1+\exp(-a)} \quad (4.59)$$

「回帰」と言っているが分類のためのモデルである

→パラメータ $\mathbf{w}$ を求める問題を解いていくものの、そうして得られた $p(C_1|\phi)$ は、クラス分類問題で利用される。

識別モデルの場合

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (4.87)$$

$M$ 個

生成モデルの場合(4.2.1)

$$p(C_1|\phi) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) \quad (4.87)$$

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad (4.66)$$

$$w_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)} \quad (4.67)$$

$M$ 個

$\frac{M(M+1)}{2}$ 個

$M$ 個

1個

$$p(C_2) = 1 - p(C_1)$$

次元が大きい場合は識別モデルの方が有利

## 4.3.2 ロジスティック回帰

ロジスティック回帰モデルのパラメータを最尤法を用いて決定する。  
そのための準備として、ロジスティックシグモイド関数の微分を求める。

(演習4.12)  $\frac{d\sigma}{da} = \sigma(1 - \sigma)$  の証明

$$\sigma = \frac{1}{1 + \exp(-a)}$$

対数を取る

$$\ln \sigma = -\ln(1 + \exp(-a))$$

合成関数の微分

$$\frac{1}{\sigma} \frac{d\sigma}{da} = \frac{\exp(-a)}{1 + \exp(-a)} = 1 - \sigma$$

対数の微分

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

両辺 $\sigma$ 倍

## 4.3.2 ロジスティック回帰

データ集合:

$$\{\phi_n, t_n\}$$

$\phi_n$ が、ロジスティック回帰モデルを用いて  $t_n$  に分類された、という訓練データ

ただし、

$$\phi_n = \phi(x_n)$$

$$t_n \in \{0, 1\}$$

$$n = 1, \dots, N$$

$t_n$ のそれぞれの尤度関数は、

$$p(t_n|\mathbf{w}) = \begin{cases} p(c_1|\phi_n) = y_n & (t_n = 1 \text{ のとき}) \\ p(c_2|\phi_n) = 1 - y_n & (t_n = 0 \text{ のとき}) \end{cases}$$

と書ける。これはさらに、 $t_n$ の離散性を用いて、以下のように表せる。

$$p(t_n|\mathbf{w}) = y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

以上より、 $N$ 個のデータに対する尤度関数は、

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N p(t_n|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad (4.89)$$

誤差関数は、負の対数を取り、

$$E(\mathbf{w}) = -\ln(p(\mathbf{t}|\mathbf{w})) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (4.90)$$

となる。これを交差エントロピー誤差関数と呼ぶ。



### 4.3.2 ロジスティック回帰

$$y_n = \sigma(a_n)$$
$$a_n = \mathbf{w}^T \phi_n$$

より、

$$\begin{aligned}\nabla \ln y_n &= \frac{d \ln y_n}{d \mathbf{w}} = \frac{d \ln y_n}{dy_n} \frac{dy_n}{da_n} \frac{da_n}{d \mathbf{w}} \\ &= \frac{1}{y_n} y_n(1 - y_n) \phi_n \\ &= (1 - y_n) \phi_n\end{aligned}$$

合成関数の微分

$$\nabla \ln(1 - y_n) = -y_n \phi_n$$

これを使って、 $E(\mathbf{w})$ の勾配は、

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (4.90)$$

$$\begin{aligned}\nabla E(\mathbf{w}) &= - \sum_{n=1}^N \{t_n(1 - y_n) \phi_n - (1 - t_n) y_n \phi_n\} \\ &= \sum_{n=1}^N (y_n - t_n) \phi_n\end{aligned} \quad (4.91)$$

と表せる(演習4.13)。

## 4.3.2 ロジスティック回帰

(3.13)との比較

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n \quad (4.91)$$

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T \quad (3.13)$$

$E(\mathbf{w})$ にだけ着目し、 $\phi_n = \phi(\mathbf{x}_n)$ 、 $y_n = \mathbf{w}^T \phi_n$ とすると、

$$\begin{aligned} \nabla E(\mathbf{w}) &= - \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n) \\ &= \sum_{n=1}^N (y_n - t_n) \phi_n \end{aligned}$$

となる。これは(4.91)と同じ形をしており、 $y_n$ の中身だけが違う。  
この違いは第4章の冒頭でしれっと登場した活性化関数((4.3)式参照)の  
違いであり、決定面を求める際には結局 $y_n = \text{定数}$ だから決定面は線形と  
なり、以後は(3.13)と本質的には同じように議論ができる。

## 4.3.2 ロジスティック回帰

### (演習4.14) 最尤解の考察

データ集合が線形分離可能であるとは、ある $\mathbf{w}_0$ があって、超平面 $\mathbf{w}_0^T \phi = 0$ が空間を分離し、 $\mathbf{w}_0^T \phi > 0$ の時 $C_1$ 、 $\mathbf{w}_0^T \phi < 0$ の時 $C_2$ という風に分類できることを言う。

このような $\mathbf{w}_0$ に対し、

$$\sigma(\mathbf{w}_0^T \phi) = \frac{1}{1 + \exp(-\mathbf{w}_0^T \phi)} \begin{cases} > 0.5 & (\phi \in C_1) \\ < 0.5 & (\phi \in C_2) \end{cases}$$

であるから、

$$\lim_{\|\mathbf{w}_0\| \rightarrow \infty} \exp(-\mathbf{w}_0^T \phi) = \begin{cases} 0 & (\phi \in C_1) \\ \infty & (\phi \in C_2) \end{cases}$$

したがって、

$$\lim_{\|\mathbf{w}_0\| \rightarrow \infty} y_n = \lim_{\|\mathbf{w}_0\| \rightarrow \infty} \frac{1}{1 + \exp(-\mathbf{w}_0^T \phi_n)} = \begin{cases} 1 & (\phi \in C_1) \\ 0 & (\phi \in C_2) \end{cases} = t_n$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n \tag{4.91}$$

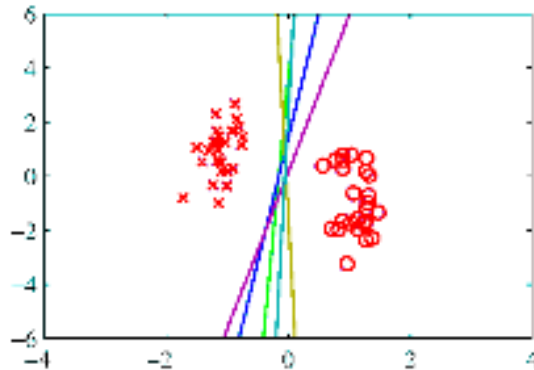
$$\therefore \lim_{\|\mathbf{w}_0\| \rightarrow \infty} \nabla E(\mathbf{w}_0) = 0$$

以上より、このような $\mathbf{w}_0$ が誤差関数の最尤解になっていることが示された。

## 4.3.2 ロジスティック回帰

### 最尤解の考察続き

データ集合が線形分離可能である時、 $\mathbf{w}^T \phi = 0$ が決定境界となっているような $\mathbf{w}$ は無数に存在する(下図)。



最尤法では、解の中の1つを他の解と区別して好んで選別することはできない。→最適化アルゴリズムと初期値に依存し、場合によっては $\mathbf{w}$ の大きさが発散してしまう。

## 4.3.2 ロジスティック回帰

### 最尤解の考察続き

p.205「データ数がパラメータ数と比べて大きくても、学習データ集合が線形分離可能である限りこの問題は発生してしまう」について

データ数がパラメータ数と比べて大きい場合、 $\phi_1, \dots, \phi_n$  はベクトルの個数 (= データ数) が空間の次元 (= パラメータ数) より大きいため、線形独立ではない。つまり、

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = 0$$

を満たし、かつ、いくつかの  $n$  については

$$y_n - t_n \neq 0$$

であるような  $\mathbf{w}$  が存在するかもしれない。

しかし、そのような場合でも、線形分離可能な場合は、アルゴリズムの選び方によって  $\mathbf{w}$  の大きさが  $\infty$  であるようなものを選んでしまう可能性がある。

1.1 や 1.2.5 と同様に罰則項 ( $\mathbf{w}$  が大きいほど  $E(\mathbf{w})$  が大きくなるようにする) を加え、正則化することで、この問題を避けることができる。

### 4.3.3 反復重み付け最小二乗

#### 最小二乗法

誤差関数 $E(w)$ の値を最小にするようなパラメータ $w$ を求める手法

#### 重み付け最小二乗法

データごとに重要度が設定されている場合の最小二乗法

#### 反復重み付け最小二乗法

データごとに重要度が設定され、それがパラメータ $w$ に依存する場合の最小二乗法。この場合、重み付け行列とパラメータとを繰り返し更新しながら最尤解を求めていく。

**線形回帰モデル**では、(ガウス分布雑音モデル、つまり雑音がガウス分布に従って現れるという仮定の下で)最尤解を解析的に導出できた。

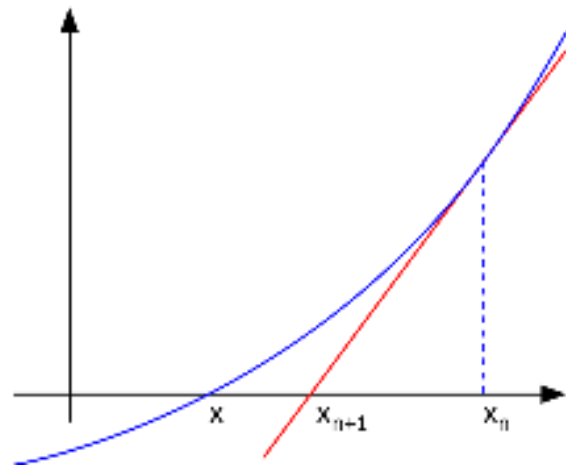
**ロジスティック回帰モデル**では、ロジスティックシグモイド関数の非線形性のために最尤解を解析的に導出できない。

→ニュートン-ラフソン法を使って近似値を求めることは可能

### 4.3.3 反復重み付け最小二乗

ニュートン-ラフソン法 (Newton-Raphson method)

入力値  $x_n$  における接線の零点を次の入力値  $x_{n+1}$  とし、これを繰り返すことで関数自体の零点を求める手法



接線の方程式:

$$y - f(x_n) = f'(x_n)(x - x_n)$$

に  $(x, y) = (x_{n+1}, 0)$  を代入し、

$$0 - f(x_n) = f'(x_n)(x_{n+1} - x_n)$$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

を漸化式として値を目標値  $x$  に近づけていく。

関数が  $M$  次元ベクトルから  $M$  次元ベクトルへの写像の場合も同様にして

$$\mathbf{0} - f(\mathbf{x}_n) = \partial f(\mathbf{x}_n)(\mathbf{x}_{n+1} - \mathbf{x}_n)$$

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \partial f(\mathbf{x}_n)^{-1} f(\mathbf{x}_n)$$

となる。 $\partial f$  はヤコビ行列と呼ばれ、以下で与えられる。

$$\partial f_{ij} = \frac{df(x_i)}{dx_j}$$

誤差関数の場合、 $f = \nabla E(\mathbf{w})$  であり、 $\partial f = \mathbf{H} = \nabla \nabla E(\mathbf{w})$  である。

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w}^{(\text{old})}) \quad (4.93)$$

$\mathbf{H}$  は  $E$  のヘッセ行列と呼ばれ、 $\mathbf{H} = \nabla \nabla E(\mathbf{w})$  と定義される。

### 4.3.3 反復重み付け最小二乗

最小二乗法の場合

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \quad (3.12)$$

この勾配とヘッセ行列は

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t} \quad (4.93)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \Phi^T \Phi \quad (4.94)$$

で与えられる。これを式(4.92)に代入し、

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \Phi)^{-1} \{ \Phi^T \Phi \mathbf{w}^{(\text{old})} - \Phi^T \mathbf{t} \} \\ &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \end{aligned} \quad (4.95)$$

この式は、 $\nabla E(\mathbf{w}) = 0$ とした時の解と一致する。



### 4.3.3 反復重み付け最小二乗

#### 重み付き最小二乗法の場合

(教科書には載っていないが)データ点 $t_n$ に重み $r_n > 0$ が割り当てられている場合について考察する。この時の誤差関数は

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N r_n \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \quad (3.104)$$

で与えられる。この勾配とヘッセ行列は

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N r_n (\mathbf{w}^T \phi_n - t_n) \phi_n = \Phi^T \mathbf{R} \Phi \mathbf{w} - \Phi^T \mathbf{R} \mathbf{t} \quad (4.93)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \Phi^T \mathbf{R} \Phi \quad (4.94)$$

で与えられる。ただし、 $\mathbf{R}$ は $R_{nn} = r_n$ を満たす対角行列である。これを式(4.92)に代入し、

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\text{old})} - \Phi^T \mathbf{R} \mathbf{t} \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{t} \end{aligned} \quad (4.95)$$

この場合も、 $\nabla E(\mathbf{w}) = 0$ とした時の解と一致する。

### 4.3.3 反復重み付け最小二乗

重み付き最小二乗法の場合： $r_n$ の解釈

誤差が正規分布に従い、その分散がデータに依存すると仮定した場合、その尤度関数は、

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \boldsymbol{\beta}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta_n^{-1})$$

負の対数を取り、 $\mathbf{w}$ に関する部分だけ取ると、

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \beta_n \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

重み付き最小二乗法の誤差関数と比較すると、結局、

$$r_n = \beta_n$$

であり、重み係数 $r_n$ は誤差の精度パラメータであると解釈できる。

oO(後の議論で「IRLSの対角行列の要素は、重み付き最小二乗法の場合と同様分散であると解釈できる」と出てくるが、これは間違いな気がしている。)

### 4.3.3 反復重み付け最小二乗

反復重み付き最小二乗法の場合

ロジスティック回帰における交差エントロピー誤差関数:

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (4.90)$$

の勾配とヘッセ行列は

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t}) \quad (4.96)$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi \quad (4.97)$$

で与えられる。ただし、 $\mathbf{R}$ は $\mathbf{R}_{nn} = y_n(1 - y_n)$ を満たす対角行列である。  
これを式(4.92)に代入し、

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\text{old})} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \{ \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}) \} \end{aligned} \quad (4.99)$$

ここで、

### 4.3.3 反復重み付け最小二乗

反復重み付き最小二乗法の場合

$$\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1}(\mathbf{y} - \mathbf{t}) \quad (4.100)$$

とおくと、

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \{ \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1}(\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z} \end{aligned} \quad (4.99)$$

さて、この式は、重み付き最小二乗法の場合の解:

$$\mathbf{w}^{(\text{new})} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{t}$$

とよく似ている。違うのは、 $\mathbf{R}$ が $\mathbf{w}$ の関数になっている点、 $\mathbf{z}$ がなんだかごちゃごちゃしている点である。

$\mathbf{R}$ が $\mathbf{w}$ の関数になっている点

このために、 $\mathbf{w}$ が新たに求まるたびに $\mathbf{R}$ を再計算し直して、正規方程式を繰り返し解かなければならない。反復重み付き最小二乗法(IRLS)の名はここからきている。

各 $t_n$ の平均、分散は、

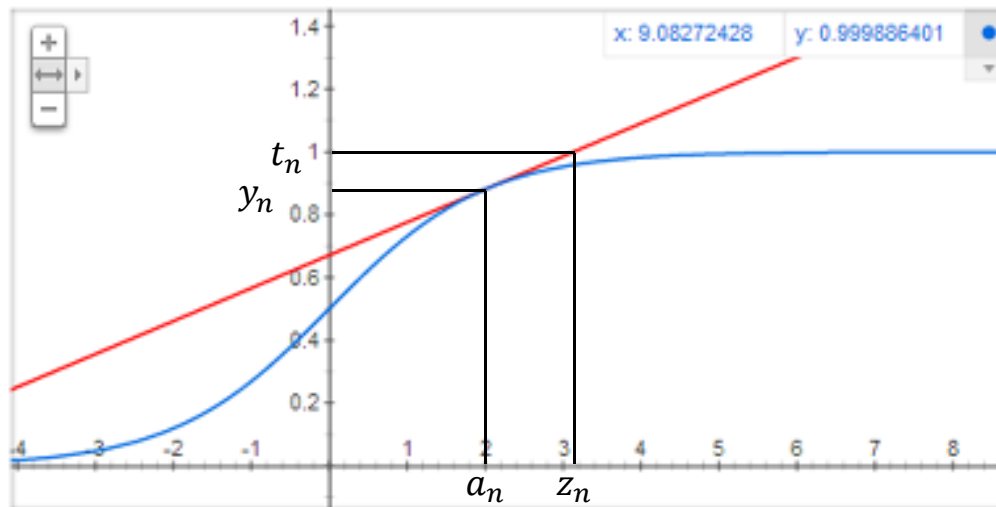
$$\mathbb{E}[t_n] = \sigma(\mathbf{w}^T \phi_n) = y_n \quad (4.101)$$

$$\text{var}[t_n] = \mathbb{E}[t_n^2] - \mathbb{E}[t_n]^2 = \sigma(\mathbf{w}^T \phi_n) - \sigma(\mathbf{w}^T \phi_n)^2 = y_n(1 - y_n) \quad (4.102)$$

だから $\mathbf{R}$ の要素は分散だと解釈できる。

### 4.3.3 反復重み付け最小二乗

zがなんだかごちゃごちゃしている点



$$y = \sigma(a) = \frac{1}{1 + \exp(-a)}$$

$z_n$ と $a_n$ 、 $y_n$ 、 $t_n$ は、上図のような位置関係にある。

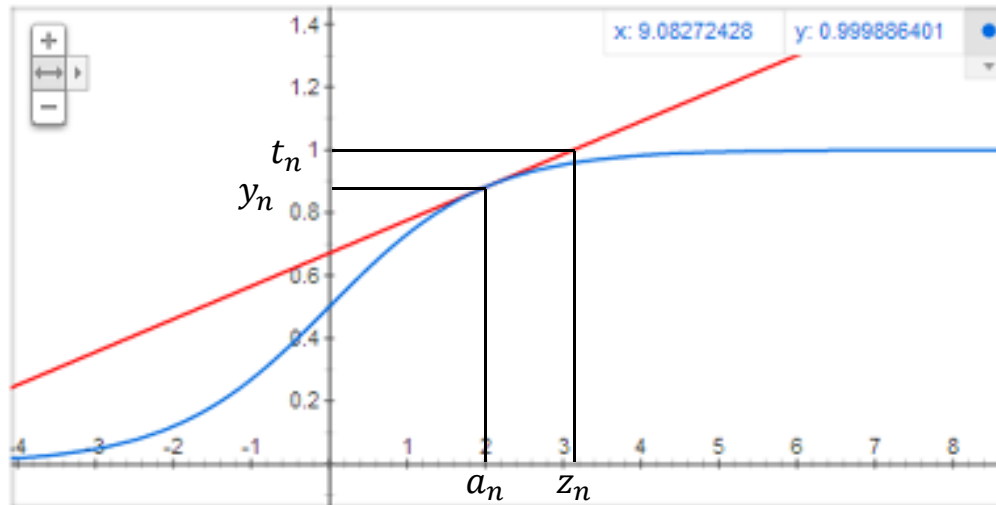
$$t_n - y_n = \left. \frac{dy}{da} \right|_{a=a_n} (z_n - a_n)$$

$\left. \frac{dy}{da} \right|_{a=a_n} = y_n(1 - y_n) = r_n$ 、 $a_n = \phi_n^T \mathbf{w}^{(\text{old})}$ を使って上式を $z_n$ について解くと、

$$\begin{aligned} t_n - y_n &= r_n(z_n - a_n) \\ z_n &= a_n + r_n^{-1}(t_n - y_n) \\ &= \phi_n^T \mathbf{w}^{(\text{old})} - r_n^{-1}(y_n - t_n) \end{aligned}$$

### 4.3.3 反復重み付け最小二乗

$\mathbf{z}$ がなんだかごちゃごちゃしている点



$$y = \sigma(a) = \frac{1}{1 + \exp(-a)}$$

$n$ についてまとめると、

$$\mathbf{z} = \Phi \mathbf{w} - \mathbf{R}^{-1}(\mathbf{y} - \mathbf{t})$$

となって、確かに(4.100)と一致する

線形回帰モデルの場合は、 $y_n$ が $\mathbf{w}$ の線形関数だったため、そのまま

$$\Phi^T \mathbf{R}(\mathbf{y} - \mathbf{t}) = \mathbf{0}$$

としてこれを解くことができた。ロジスティック回帰モデルの場合、ロジスティックシグモイド関数を線形近似し、そこでの解を求めると解釈できる。

$$\Phi^T \mathbf{R}(\mathbf{a}^{(\text{new})} - \mathbf{z}) = \Phi^T \mathbf{R}(\Phi \mathbf{w}^{(\text{new})} - \mathbf{z}) = \mathbf{0}$$

$$\mathbf{w}^{(\text{new})} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}$$

## 4.3.4 多クラスロジスティック回帰

### ソフトマックス変換

多クラスの場合、事後確率は特徴変数の線形関数のソフトマックス変換で与えられる:

$$p(C_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad (4.104)$$

ただし、

事後確率は

ソフトマックス変換で与えられる

$$a_k = \mathbf{w}_k^T \phi \quad (4.105)$$

特徴変数の線形関数の

2クラス分類と同じように、ベイズ定理を用いる生成モデルに対し、最尤法を用いてパラメータ $\{\mathbf{w}_k\}$ を決定する。

### 4.3.4 多クラスロジスティック回帰

準備:(演習4.17)  $\frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$  の証明

(4.104)の両辺対数を取り、 $a_j$ で微分する(分母の総和変数を*i*とした)

$$\ln y_k = a_k - \ln \sum_i \exp a_i$$
$$\frac{1}{y_k} \frac{\partial y_k}{\partial a_j} = \frac{\partial a_k}{\partial a_j} - \frac{\exp a_j}{\sum_i \exp a_i}$$

$a_j$ による偏微分のため  
添え字*j*のみが残る

ここで、

$$\frac{\partial a_k}{\partial a_j} = \begin{cases} 1 & (j = k) \\ 0 & (j \neq k) \end{cases} = I_{kj}$$

合成関数の微分

より、

$$\frac{1}{y_k} \frac{\partial y_k}{\partial a_j} = I_{kj} - y_j$$
$$\therefore \frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$$



## 4.3.4 多クラスロジスティック回帰

尤度関数の表現: 1-of-K符号化法のうまみ

目的変数 $t_n$ は $K$ 次元ベクトルで、データ $\mathbf{x}_n$ がクラス $k_0$ に分類される場合、

$$t_{nk} = \begin{cases} 1 & (k = k_0) \\ 0 & (k \neq k_0) \end{cases}$$

という形をしている。これを使うと、尤度関数は簡単に表すことができ、

$$p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k|\phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \quad (4.107)$$

$t_{nk}$ を要素とする $N \times K$ 行列

たくさんかけているように見えるが  
実は $K$ 個のうち1つしか有効でない

$y_{nk} = y_k(\phi_n)$

誤差関数は、

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T}|\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \quad (4.108)$$

と表せる。これを(多クラス分類問題の)交差エントロピー関数と呼ぶ。

## 4.3.4 多クラスロジスティック回帰

演習4.18の証明

$$\begin{aligned}\nabla_{\mathbf{w}_j} \ln y_{nk} &= \frac{d \ln y_{nk}}{dy_{nk}} \frac{\partial y_{nk}}{\partial a_{nj}} \frac{da_{nj}}{d\mathbf{w}_j} \\ &= \frac{1}{y_{nk}} y_{nk} (I_{kj} - y_{nj}) \phi_n = (I_{kj} - y_{nj}) \phi_n\end{aligned}$$

$$\sum_{k=1}^K t_{nk} I_{kj} = t_{nj} I_{jj} = t_{nj}$$

$$\sum_{k=1}^K t_{nk} y_{nj} = y_{nj} \sum_{k=1}^K t_{nk} = y_{nj} \cdot \mathbf{1} = y_{nj}$$

よって、

$$\begin{aligned}\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) &= \nabla_{\mathbf{w}_j} \left( - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \right) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \nabla_{\mathbf{w}_j} \ln y_{nk} \\ &= - \sum_{n=1}^N \sum_{k=1}^K t_{nk} (I_{kj} - y_{nj}) \phi_n \\ &= \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n\end{aligned}$$

## 4.3.4 多クラスロジスティック回帰

### ニュートン-ラフソン法の適用

多クラスの場合、パラメータ数は、 $M$ 次元ベクトル $\mathbf{w}_k$ が $K$ 個で $MK$ 個ある。  
これらを縦に1列に並べて、

$$\mathbf{v} = (\mathbf{w}_1^T, \dots, \mathbf{w}_K^T)^T$$

と置くと、2クラスの場合と同様にニュートン-ラフソン法を適用できる。

$$\mathbf{v}^{(\text{new})} = \mathbf{v}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{v})$$

$\mathbf{H}$ は $MK \times MK$ 行列で、次のようにあらわされる。

$$\mathbf{H} = \begin{pmatrix} \mathbf{h}_{11} & \cdots & \mathbf{h}_{K1} \\ \vdots & \ddots & \vdots \\ \mathbf{h}_{1K} & \cdots & \mathbf{h}_{KK} \end{pmatrix}$$

$$\mathbf{h}_{jk} = \nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N y_{nk} (I_{kj} - y_{nk}) \phi_n \phi_n^T$$

### 4.3.4 多クラスロジスティック回帰

演習4.20の証明

$$(\mathbf{R}_{jk})_{nn} = y_{nk}(I_{kj} - y_{nj})$$

と置くと、

$$\mathbf{h}_{jk} = \mathbf{\Phi}^T \mathbf{R}_{jk} \mathbf{\Phi}$$

と書ける。

$$\mathbf{R} = \begin{pmatrix} \mathbf{R}_{11} & \cdots & \mathbf{R}_{K1} \\ \vdots & \ddots & \vdots \\ \mathbf{R}_{1K} & \cdots & \mathbf{R}_{KK} \end{pmatrix}, \quad \mathbf{\Psi} = \begin{pmatrix} \mathbf{\Phi} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{\Phi} \end{pmatrix}$$

と置くと、 $\mathbf{R}$ は $NK \times NK$ 行列、 $\mathbf{\Psi}$ は $NK \times MK$ 行列になり、

$$\mathbf{H} = \mathbf{\Psi}^T \mathbf{R} \mathbf{\Psi}$$

と書ける。これが正定値であることを証明する。まず、 $j \neq k$ の時

$$(\mathbf{R}_{jk})_{nn} = -y_{nk}y_{nj} = (\mathbf{R}_{kj})_{nn}$$

より、 $\mathbf{R}$ は対称行列である。

$MK$ 次元ベクトル $\mathbf{u}$ に対し、

$$\mathbf{v} = \mathbf{\Psi} \mathbf{u} = (\mathbf{v}_1^T, \dots, \mathbf{v}_K^T)^T$$

と置くと、

$$\mathbf{u}^T \mathbf{H} \mathbf{u} = \mathbf{v}^T \mathbf{R} \mathbf{v} = (\mathbf{v}_1^T, \dots, \mathbf{v}_K^T)^T \begin{pmatrix} \mathbf{R}_{11} & \cdots & \mathbf{R}_{K1} \\ \vdots & \ddots & \vdots \\ \mathbf{R}_{1K} & \cdots & \mathbf{R}_{KK} \end{pmatrix} \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_K \end{pmatrix}$$

## 4.3.4 多クラスロジスティック回帰

演習4.20の証明

$$\begin{aligned}\mathbf{u}^T \mathbf{H} \mathbf{u} &= \mathbf{v}^T \mathbf{R} \mathbf{v} = (\mathbf{v}_1^T, \dots, \mathbf{v}_K^T)^T \begin{pmatrix} \mathbf{R}_{11} & \cdots & \mathbf{R}_{K1} \\ \vdots & \ddots & \vdots \\ \mathbf{R}_{1K} & \cdots & \mathbf{R}_{KK} \end{pmatrix} \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_K \end{pmatrix} \\&= (\mathbf{v}_1^T, \dots, \mathbf{v}_K^T)^T \begin{pmatrix} \mathbf{R}_{11} \mathbf{v}_1 + \cdots + \mathbf{R}_{K1} \mathbf{v}_K \\ \vdots \\ \mathbf{R}_{K1} \mathbf{v}_1 + \cdots + \mathbf{R}_{KK} \mathbf{v}_K \end{pmatrix} \\&= \sum_{j=1}^K \sum_{k=1}^K \mathbf{v}_j^T \mathbf{R}_{jk} \mathbf{v}_k \\&= \sum_{j=1}^K \sum_{k=1}^K \sum_{n=1}^N y_{nk} (I_{jk} - y_{nj}) v_{jn} v_{kn} \\&= \sum_{n=1}^N \sum_{j=1}^K \sum_{k=1}^K y_{nk} (I_{jk} - y_{nj}) v_{jn} v_{kn} \\&= \sum_{n=1}^N \sum_{j=1}^K \sum_{k=1}^K (y_{nk} I_{jk} v_{jn} v_{kn} - y_{nk} y_{nj} v_{jn} v_{kn})\end{aligned}$$

## 4.3.4 多クラスロジスティック回帰

演習4.20の証明

$$\sum_{j=1}^K \sum_{k=1}^K y_{nk} I_{jk} v_{jn} v_{kn} = \sum_{j=1}^K y_{nj} v_{jn} v_{jn} = \sum_{j=1}^K \sum_{k=1}^K y_{nk} y_{nj} v_{jn} v_{jn}$$

より、

$$\begin{aligned} \mathbf{u}^T \mathbf{H} \mathbf{u} &= \sum_{n=1}^N \sum_{j=1}^K \sum_{k=1}^K (y_{nk} I_{jk} v_{jn} v_{kn} - y_{nk} y_{nj} v_{jn} v_{kn}) \\ &= \sum_{n=1}^N \sum_{j=1}^K \sum_{k=1}^K (y_{nk} y_{nj} v_{jn} v_{jn} - y_{nk} y_{nj} v_{jn} v_{kn}) \\ &= \sum_{n=1}^N \sum_{j=1}^K \sum_{k=1}^K y_{nk} y_{nj} v_{jn} (v_{jn} - v_{kn}) \end{aligned}$$

$a_{jk} = y_{nk} y_{nj} v_{jn} (v_{jn} - v_{kn})$ とおくと、

$$a_{jj} = 0,$$

$$a_{jk} + a_{kj} = y_{nk} y_{nj} v_{jn} (v_{jn} - v_{kn}) + y_{nj} y_{nk} v_{kn} (v_{kn} - v_{jn})$$

$$= y_{nk} y_{nj} (v_{jn} - v_{kn})^2$$

以上により、

### 4.3.4 多クラスロジスティック回帰

演習4.20の証明

$$\begin{aligned}\mathbf{u}^T \mathbf{H} \mathbf{u} &= \sum_{n=1}^N \sum_{j=1}^K \sum_{k=1}^K y_{nk} y_{nj} v_{jn} (v_{jn} - v_{kn}) \\ &= \sum_{n=1}^N \sum_{1 \leq j < k \leq K} y_{nk} y_{nj} (v_{jn} - v_{kn})^2 > 0\end{aligned}$$

したがって、 $H$ は正定値行列であり、誤差関数は唯一の最小解を持つ。

## 4.3.5 プロビット回帰

もう少し一般化した仮定のもとで議論する

活性化関数をもう少し一般的に $f(a)$ とし、 $a = \mathbf{w}^T \phi$ として

$$p(t = 1|a) = f(a)$$

と書ける場合を考える。例として、雑音しきい値モデルについて考察する。

雑音しきい値モデル

各入力 $\phi_n$ に対して $a_n = \mathbf{w}^T \phi_n$ を評価し、あるしきい値 $\theta$ との大小で目標編数値を設定する。

$$\begin{cases} t_n = 1 & a_n \geq \theta \text{ のとき} \\ t_n = 0 & a_n < \theta \text{ のとき} \end{cases}$$

$\theta$ がある確率密度 $p(\theta)$ から得られると仮定すると、活性化関数は

$$f(a) = p(t = 1|a) = p(a \geq \theta) = \int_{-\infty}^a p(\theta) d\theta$$

と累積分布関数の形に書ける。

プロビット関数

$p(\theta)$ の例として、標準正規分布を考える

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta|0,1) d\theta$$

これの逆関数をプロビット関数と言う。形は図4.9を参照



## 4.3.5 プロビット回帰

### 演習4.21: erf関数との関係

以下の関数をerf関数または誤差関数という。

$$\operatorname{erf}(a) = \frac{1}{\sqrt{\pi}} \int_0^a \exp(-\theta^2) d\theta$$

プロビット関数の逆関数との関係は、

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(\theta|0,1) d\theta = \left( \int_{-\infty}^0 + \int_0^a \right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\theta^2}{2}\right) d\theta$$

右辺第1項は $\frac{1}{2}$ 、第2項は $\phi = \frac{\theta}{\sqrt{2}}$ と変数変換し、

$$\begin{aligned} \Phi(a) &= \frac{1}{2} + \int_0^{\frac{a}{\sqrt{2}}} \frac{1}{\sqrt{\pi}} \exp(-\phi^2) d\theta \\ &= \frac{1}{2} + \frac{1}{2} \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \\ &= \frac{1}{2} \left\{ 1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right\} \end{aligned}$$

## 4.3.5 プロビット回帰

外れ値について

$$\sigma'(a) = \frac{\exp(-a)}{(1 + \exp(-a))^2} \cong \exp(-a)$$

$$\Phi'(a) = \mathcal{N}(a|0,1) \cong \exp(-a^2)$$

より、プロビット関数の方が外れ値に敏感である。

たとえば、 $a_n = \mathbf{w}^T \phi_n = \sqrt{2}$ ,  $t_n = 0$ となるデータ点に対しては、

$$\sigma(a_n) \doteq 0.804$$

$$\Phi(a_n) = \frac{1}{2} \{1 + \operatorname{erf}(1)\} \doteq 0.921$$

(google電卓機能を利用)

となって、プロビット関数の方が外れ値として大きく加算されてしまう。

目的変数値 $t$ が間違った値に反転する確率を $\epsilon$ とすると、データ点 $\mathbf{x}$ における目的変数値の確率分布は

$$\begin{aligned} p(t = 1|\mathbf{x}) &= (1 - \epsilon)\sigma(\mathbf{x}) + \epsilon(1 - \sigma(\mathbf{x})) \\ &= \epsilon + (1 - 2\epsilon)\sigma(\mathbf{x}) \end{aligned}$$

ただし、 $\epsilon$ は事前に設定されるか、データ点から推定されるハイパーパラメータとして扱われる

## 4.3.6 正準連結関数

より一般的な議論

これまでに取り扱ったモデルはすべて、誤差関数の微分が同じ形式:

$$\sum_{n=1}^N \{y_n - t_n\} \phi_n$$

を取った。

この節では、上記の結論が、次の仮定のもとに得られることを示す

(仮定1) 目的変数に対する条件付き確率分布を指数型分布族の中から選ぶ。

(仮定2) 期待値 $y$ と線形予測子 $\mathbf{w}^T \phi$ との連結関数に、正準連結関数を選ぶ。

線形予測子 (linear predictor)

独立変数の一次式 ( $\mathbf{w}^T \phi$ )

連結関数

$\mathbf{w}^T \phi$ を、目的変数 $t$ の期待値 $y$ の関数で表したもの (活性化関数  
 $y = f(\mathbf{w}^T \phi)$ の逆関数)

正準連結関数

連結関数の中で特に、目的変数の分布を表すモデルのパラメータ $\eta$ と期待値 $y$ との関数になっているもの

## 4.3.6 正準連結関数

より一般的な議論

仮定1より、目的変数 $t$ の条件付き確率分布は以下で与えられる。

$$p(t|\eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\} \quad (4.118)$$

ただし、 $s$ は尺度パラメータで、全ての観測点で共通だと仮定する。  
2.4と同様の議論をして $t$ の期待値を求めたい。(2.194),(2.226):

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})\} \quad (2.194)$$

$$-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}[\mathbf{u}(\mathbf{x})] \quad (2.226)$$

と比較して、

$$\begin{aligned} \mathbb{E}\left[\frac{t}{s} \middle| \eta\right] &= -\frac{d}{d\eta} \ln g(\eta) \\ y \equiv \mathbb{E}[t|\eta] &= -s \frac{d}{d\eta} \ln g(\eta) \end{aligned} \quad (4.119)$$

を得る。 $y$ が $\eta$ の関数になっている。この逆関数が存在するとしてそれを $\eta = \psi(y)$ とする。

これの対数尤度関数は、

$$\ln p(t|\eta, s) = \sum_{n=1}^N \ln p(t_n | \eta, s) = \sum_{n=1}^N \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} + \text{const.}$$

と書ける。

## 4.3.6 正準連結関数

より一般的な議論

一般化線形モデルの場合、目的変数 $t$ の期待値 $y$ が、特徴ベクトルの一次結合の非線形変換で与えられるとしている。つまり、

$$y = f(\mathbf{w}^T \phi)$$

と書ける。この場合、 $y$ の勾配は、 $a = \mathbf{w}^T \phi$

$$\nabla_{\mathbf{w}} y = \frac{dy}{da} \frac{da}{d\mathbf{w}} = f'(a) \phi$$

と書ける。これを使って、対数尤度関数の $\mathbf{w}$ に関する微分を求めると、

$$\begin{aligned} \nabla_{\mathbf{w}} \ln p(\mathbf{t}|\boldsymbol{\eta}, s) &= \frac{d}{d\mathbf{w}} \left( \sum_{n=1}^N \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} \right) \\ &= \sum_{n=1}^N \left( \frac{d}{d\eta_n} \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} \right) \frac{d\eta_n}{dy_n} \frac{dy_n}{da_n} \frac{da_n}{d\mathbf{w}} \\ &= \sum_{n=1}^N \left\{ -\frac{y_n}{s} + \frac{t_n}{s} \right\} \psi'(y_n) f'(a_n) \phi_n \end{aligned}$$

$y_n = -s \frac{d}{d\eta_n} \ln g(\eta)$

$\eta = \psi(y)$

$y = f(a_n)$

## 4.3.6 正準連結関数

より一般的な議論

さて、連結関数 $f^{-1}$ は、一般化線形モデルの上では何でもよいとしてきたが、ここで、正準連結関数を選ぶ、と仮定する。つまり、

$$f^{-1}(y) = \psi(y) \quad (4.121)$$

この仮定はとても強力で、あとの議論がずっと楽になる。

$$y = f(\psi(y))$$

の両辺を微分して、

$$1 = f'(\psi(y))\psi'(y)$$

$a = f^{-1}(y) = \psi$ より、

$$f'(a)\psi'(a) = 1$$

これを誤差関数の勾配に代入して、

$$\begin{aligned} \nabla_{\mathbf{w}} E(\mathbf{w}) &= -\nabla_{\mathbf{w}} \ln p(t|\eta, s) = -\sum_{n=1}^N \left\{ -\frac{y_n}{s} + \frac{t_n}{s} \right\} \psi'(y_n) f'(a_n) \phi_n \\ &= \frac{1}{s} \sum_{n=1}^N \{y_n - t_n\} \phi_n \end{aligned}$$

以上により、指数型分布族とそれに対応する正準連結関数を選べば、その誤差関数の勾配が誤差 $y_n - t_n$ と特徴ベクトル $\phi_n$ との積で表せることが分かった。

## 4.3.6 正準連結関数

指数分布族と正準連結関数の対照表(以下で、 $a = \mathbf{w}^T \phi$ )

モデル名	構造	正準連結関数	活性化関数
正規分布	$\mathcal{N}(t y, \beta^{-1})$	$a = id_y = y$	$y = id_a = a$
ベルヌーイ分布	$\text{Bern}(t y) = y^t(1-y)^{1-t}$	$a = \ln \frac{y}{1-y}$	$y = \sigma(a) = \frac{1}{1 + \exp(-a)}$
多項分布	$\prod_{k=1}^K y_k^{t_k}$	$a_k - a_K = \ln \frac{y_k}{y_K}$ ( $\sum_j y_j = 1$ )	$y_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$
ポアソン分布	$\frac{\lambda^t e^{-\lambda}}{t!}$	$a = \ln y$	$y = \exp a$
指数分布	$\lambda e^{-\lambda t} \ (0 \leq t)$	$a = \frac{1}{y}$	$y = \frac{1}{a}$

例: ポアソン分布

$$\frac{\lambda^t e^{-\lambda}}{t!} = \frac{1}{t!} e^{-\lambda} \exp(t \ln \lambda)$$

指数型分布族の形と照らし合わせて、

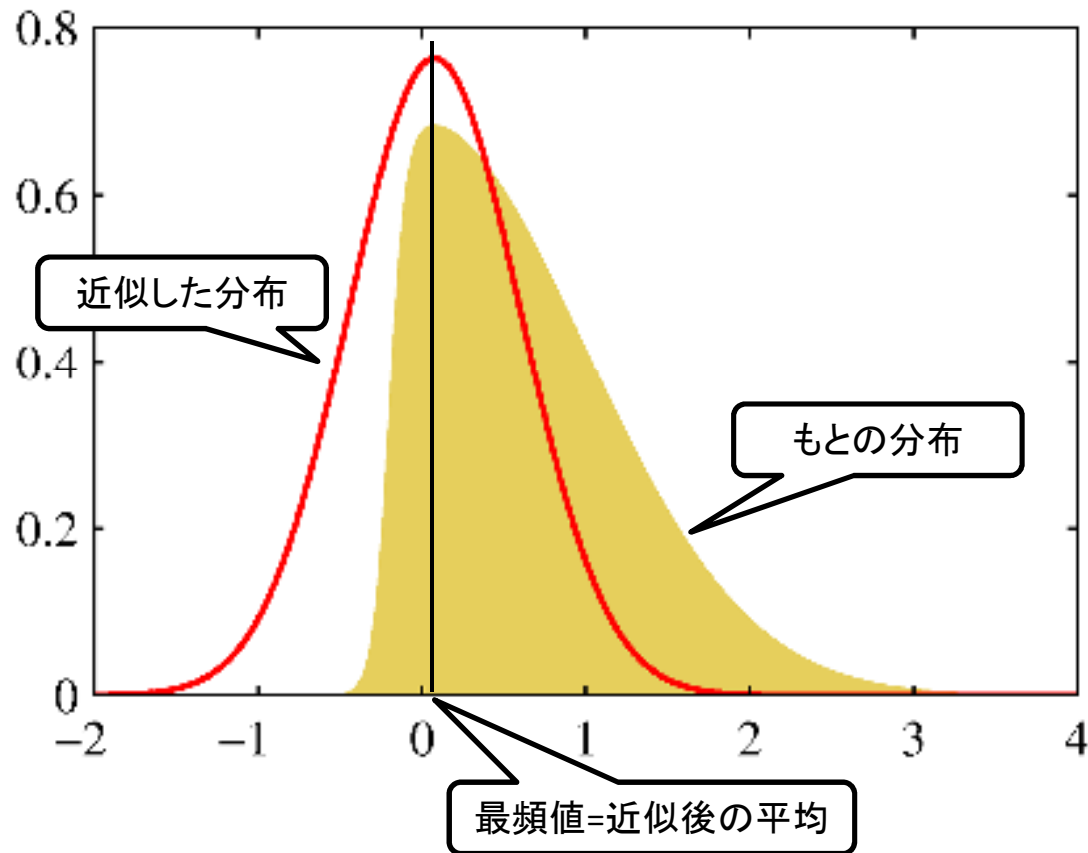
$$\eta = \ln \lambda, \quad h(t) = \frac{1}{t!}, \quad g(\eta) = e^{-\lambda} = e^{-\exp \eta}$$

$$y = -\frac{d}{d\lambda} \ln g(\eta) = -\frac{d}{d\lambda} (-\exp \eta) = \exp \eta$$

## 4.4 ラプラス近似

### ラプラス近似とは

分布を正規分布で近似してしまおうという大胆な発想  
最頻値を近似後の正規分布の平均値として、うまいこと近似をする。





## 4.4 ラプラス近似

近似のころ

- ・最頻値における1次微分は0である
  - ・正規分布の指数の肩には変数の2次関数が乗っている
- もとの分布の対数を取って2次で近似すればうまくいきそう

ごりごり計算

$$p(z) = \frac{1}{Z} f(z)$$

とする。(  $Z = \int f(z) dz$  は正規化係数、未知)  
最頻値を  $z_0$  とすると、

$$\left. \frac{df}{dz} \right|_{z=z_0} = 0$$

が成り立つ。したがって、 $\ln f$  のテイラー展開には1次の項が現れず、

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A (z - z_0)^2$$

と書ける。ここで、

$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

である。

## 4.4 ラプラス近似

ごりごり計算

両辺の指数をとると、

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$

と近似することができる。指数部分だけ取り出して、正規化した分布 $q(z)$ を作ると、これは、

$$q(z) = \left( \frac{A}{2\pi} \right)^{\frac{1}{2}} \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\}$$

と書ける。

この近似がうまくいくのは、 $A > 0$ の場合、つまり、定常点 $z_0$ が局所最大である場合である。

## 4.4 ラプラス近似

### 多次元の場合

今度は、多次元の場合：

$$p(\mathbf{z}) = \frac{f(\mathbf{z})}{Z}$$

の場合を考える。この場合も、最頻値 $\mathbf{z}_0$ のまわりでの勾配が0であること、を利用して、 $\ln f(\mathbf{z})$ を $\mathbf{z}_0$ の周りでテイラー展開すると、

$$\begin{aligned}\ln f(\mathbf{z}) &\simeq \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \\ f(\mathbf{z}) &\simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\}\end{aligned}$$

ただし、

$$\mathbf{A} = -\nabla \nabla \ln f(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0}$$

1次元の時と同様に正規化して、

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1})$$

を得る。ただし、 $|\mathbf{A}|$ は $\mathbf{A}$ の行列式、 $M$ は変数空間の次元である  
この近似がうまくいくには、 $\mathbf{A}$ は正定値行列、つまり、定常点 $\mathbf{z}_0$ が局所最大であって、局所最小であったり鞍点でないことが必要である。

## 4.4 ラプラス近似

### ラプラス近似のメリット

中心極限定理より、観測データが増えるほどガウス分布による近似精度が良くなる。

近似をする際に、真の分布の正規化係数 $Z$ を求める必要がない  
( $Z = \int f(z)dz$ は一般には計算することができない)。

### ラプラス近似のデメリット

実数変数の場合にしか適用できない(離散分布では直接適用できない)  
真の分布のある1点(最頻値 $z_0$ )のまわりしか注目していないため、多峰的な分布の特性のすべてを捉えることはできない。

## 4.4.1 モデルの比較とBIC

正規化係数 $Z$ の近似

(4.133)の両辺を $Z$ で割ると、

$$p(\mathbf{z}) = \frac{f(\mathbf{z})}{Z} \simeq \frac{f(\mathbf{z}_0)}{Z} \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\}$$

これと(4.134):

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} \quad (4.134)$$

を比較して、

$$\begin{aligned} \frac{f(\mathbf{z}_0)}{Z} &\simeq \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \\ Z &\simeq f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \end{aligned} \quad (4.135)$$

を得る。

## 4.4.1 モデルの比較とBIC

### ベイズモデル比較復習

$\mathcal{D}$ : データ集合

$\{(\mathcal{M}_i, \boldsymbol{\theta}_i)\}$ : モデルとパラメータの組

データ $\mathcal{D}$ は、モデル $\mathcal{M}_i$ のどれかに従って生成されたが、そのどれかはわからない。どのモデルが一番確からしいかは、モデルの事前分布 $p(\mathcal{M}_i)$ を仮定し、データ集合 $\mathcal{D}$ が与えられた時、

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{M}_i)p(\mathcal{D}|\mathcal{M}_i)$$

で評価される。右辺の $p(\mathcal{D}|\mathcal{M}_i)$ をモデルエビデンスという。

### モデルエビデンスのラプラス近似

モデルエビデンスは、周辺尤度とも呼ばれ、パラメータ $\boldsymbol{\theta}_i$ で尤度関数を周辺化して得られる:

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\boldsymbol{\theta}_i, \mathcal{M}_i)p(\boldsymbol{\theta}_i|\mathcal{M}_i)d\boldsymbol{\theta}_i$$

以下、簡単のため $\mathcal{M}_i$ を省略して、

$$p(\mathcal{D}) = \int p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$$

と書く。

## 4.4.1 モデルの比較とBIC

### 演習4.22

$$f(\boldsymbol{\theta}) = p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}), \quad Z = p(\mathcal{D})$$

とする。事後確率 $p(\boldsymbol{\theta}|\mathcal{D})$ は、ベイズの定理より、

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} = \frac{f(\boldsymbol{\theta})}{Z}$$

で与えられる。事後確率の最頻値を $\boldsymbol{\theta}_{MAP}$ としてラプラス近似を適用すると、(4.135)より、

$$\begin{aligned} Z &= \int f(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &\simeq f(\boldsymbol{\theta}_{MAP}) \int \exp\left\{-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})^T \mathbf{A}(\boldsymbol{\theta} - \boldsymbol{\theta}_{MAP})\right\} d\boldsymbol{\theta} \\ &= f(\boldsymbol{\theta}_{MAP}) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \end{aligned}$$

両辺の対数を取って、

$$\begin{aligned} \ln p(\mathcal{D}) = \ln Z &\simeq \ln f(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}| \\ &= \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}| \end{aligned} \quad (4.137)$$

となって、(4.137)を得る。

## 4.4.1 モデルの比較とBIC

### 演習4.22

ここで、 $\mathbf{A}$ は、(4.132)より、

$$\begin{aligned}\mathbf{A} &= -\nabla\nabla \ln f(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{MAP}} \\ &= -\nabla\nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP})p(\boldsymbol{\theta}_{MAP}) \\ &= -\nabla\nabla \ln p(\boldsymbol{\theta}_{MAP}|\mathcal{D})p(\mathcal{D}) \\ &= -\nabla\nabla \ln p(\boldsymbol{\theta}_{MAP}|\mathcal{D})\end{aligned}\tag{4.138}$$

である。

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2}\ln(2\pi) - \frac{1}{2}\ln|\mathbf{A}|\tag{4.137}$$

最適なパラメータを  
使用して評価した対数尤度

Occam係数

(4.137)で、右辺第1項は、最適なパラメータを使用して評価した対数尤度である。

残りの3項は「Occam係数」と呼ばれ、モデルの複雑さにペナルティーを科す役目を果たす。



## 4.4.1 モデルの比較とBIC

### ベイズ情報量基準(演習4.23)

対称行列 $\mathbf{A}$ が非退化であるとは、 $\mathbf{A}$ の行列式が0でないことを言う。パラメータの事前分布がガウス分布に従う:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{V}_0) = \frac{1}{(2\pi)^{M/2} |\mathbf{V}_0|^{1/2}} \exp \left\{ -\frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) \right\}$$

と仮定する。この両辺対数を取り、

$$\ln p(\boldsymbol{\theta}) = -\frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{V}_0| - \frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m})$$

2階微分は、

$$\nabla \nabla \ln p(\boldsymbol{\theta}) = -\mathbf{V}_0^{-1}$$

次に、 $\boldsymbol{\theta}_{MAP}$  で評価された負の対数尤度  $\ln p(\mathcal{D}|\boldsymbol{\theta})$  の2階微分を  $\mathbf{H}$  とする:

$$\mathbf{H} = -\nabla \nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP})$$

(4.138)より、

$$\begin{aligned} \mathbf{A} &= -\nabla \nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) p(\boldsymbol{\theta}_{MAP}) \\ &= -\nabla \nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \nabla \nabla \ln p(\boldsymbol{\theta}_{MAP}) = \mathbf{H} + \mathbf{V}_0^{-1} \end{aligned}$$

これらを式(4.137)に代入すると、

$$\begin{aligned} \ln p(\mathcal{D}) &\simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}| \\ &= \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{V}_0| \\ &\quad - \frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{H} + \mathbf{V}_0^{-1}| \end{aligned}$$

## 4.4.1 モデルの比較とBIC

ベイズ情報量基準(演習4.23)

$$\begin{aligned}\ln p(\mathcal{D}) &\simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{V}_0| \\ &\quad - \frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{H} + \mathbf{V}_0^{-1}| \\ &= \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) - \frac{1}{2} \ln |\mathbf{H} + \mathbf{V}_0^{-1}| + \text{const.}\end{aligned}$$

ここで、事前確率が広い幅を持っている、つまり $\mathbf{V}_0^{-1}$ が十分に小さいと仮定すると、右辺第3項の行列式は $|\mathbf{H}|$ と近似でき、

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) - \frac{1}{2} \ln |\mathbf{H}| + \text{const.}$$

を得る(演習4.23の前半)

さらに、各データ点が独立同分布だと仮定すると、

$$\begin{aligned}p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) &= \prod_{n=1}^N p(d_n|\boldsymbol{\theta}_{MAP}) \\ \mathbf{H} &= -\nabla \nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{MAP}) = \sum_{n=1}^N -\nabla \nabla \ln p(d_n|\boldsymbol{\theta}_{MAP})\end{aligned}$$

## 4.4.1 モデルの比較とBIC

ベイズ情報量基準(演習4.23)

$$\hat{\mathbf{H}} = \frac{1}{N} \sum_{n=1}^N -\nabla \nabla \ln p(d_n | \boldsymbol{\theta}_{MAP})$$

とおくと、

$$\mathbf{H} = N\hat{\mathbf{H}}$$

$$\ln|\mathbf{H}| = \ln|N\hat{\mathbf{H}}| = \ln N^M |\hat{\mathbf{H}}| = M \ln N + \ln|\hat{\mathbf{H}}|$$

さて、 $V_0^{-1}$ が十分小さいことから

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) - \frac{1}{2} \ln|\mathbf{H}| + \text{const.}$$

の右辺第2項は無視され、さらに、 $N$ が十分大きいと仮定すれば、 $\ln|\hat{\mathbf{H}}|$ も無視できる。定数項を省略して、

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D} | \boldsymbol{\theta}_{MAP}) - \frac{1}{2} M \ln N \quad (4.137)$$

を得る。

この右辺の値をベイズ情報量規準 (Bayesian Information Criterion, BIC) あるいは種ワルツ規準と呼ぶ。(1.73)のAIC:

$$\ln p(\mathcal{D} | \mathbf{w}_{ML}) - M$$

と比べると、BICの方が、罰金項が $\frac{N}{2}$ 倍大きく、モデルの複雑さにより重いペナルティを科していると言える。