

4.2.3 離散特徴

4.2.1、4.2.2などでは、特徴量 \mathbf{x} をガウス分布と仮定した際に、事後分布が一般化線形モデルで表されることを示したが、特徴量 \mathbf{x} が離散値の場合でも一般化線形モデルとなることを、 $x_i \in \{0, 1\}$ の例で示す。

特徴量の数がD個であるとする、 2^D 個の要素の表のすべての要素の確率を考えることになる。**ナイーブベイズ**を仮定すると、各特徴量がクラスに対して条件付き独立とすることができ、独立変数の数を減らすことができる。

$$\begin{aligned} P(\mathbf{x} | C_k) &= P(x_1 | C_k) P(x_2 | C_k) \cdots P(x_D | C_k) \quad (\because \text{条件付き独立}) \\ &= \prod_{i=1}^D P(x_i | C_k) \\ &= \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i} \end{aligned}$$

ここで、 $P(x_i | C_k)$ にはベルヌーイ分布を仮定している。
これを、4.63式に入れると、ソフトマックス関数を使ったこれまでの議論通りに適用できる。

$$\begin{aligned} Q_k &= \ln \left(\underbrace{P(\mathbf{x} | C_k)}_{\substack{\text{代入する} \\ \text{（算数）}}} P(C_k) \right) \quad (4.63) \\ &= \sum_{i=1}^D \{ x_i \ln \mu_{ki} + (1 - x_i) \ln (1 - \mu_{ki}) \} + P(C_k) \end{aligned}$$

これは、 $\mathbf{w} \cdot \mathbf{x}$ の線形関数である。

4.2.4 指数型分布族

ここまでで、ガウス分布と離散値入力において、事後確率 $p(x|c_k)$ が一般化線形モデルで表されることを示した。

本項ではより一般化して、クラスの条件付き確率 $p(x|c_k)$ が指数型分布族であるならば、事後分布が一般化線形モデルとなることに拡張する。

指数型分布族は、尺度パラメータ S を導入して次のように表すことができる。

$$\begin{aligned} p(x|\lambda_k) &= h(x) g(\lambda_k) \exp\{\lambda_k^T u(x)\} \\ &= h(x) g(\lambda_k) \exp\{\lambda_k^T x\} \end{aligned}$$

(相談)

「 $u(x) = x$ と仮定するような分布の範囲クラスに注目すると、
とあるが、条件付けて良いのか。任意の指数型分布族について
議論したかったのではなかったか。」

$$= \frac{1}{S} h\left(\frac{1}{S} x\right) g(\lambda_k) \exp\left\{\frac{1}{S} \lambda_k^T x\right\}$$

(相談)

ここで、(2.236) p116 の尺度パラメータを導入すると。
とあるが、これを用いる根拠を良く分っていない。

$$p(x|\sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$$

これを、先ほど同様に4.63式に代入すると、ソフトマックス関数を使った事後分布を構成することができる。その際の a_k は、

$$\alpha_k(x) = \ln(p(x|c_k)p(c_k)) \quad (4.63)$$

$$= \underbrace{-\ln(s) + \ln(h(\frac{x}{s}))}_{k \text{ に依存しない部分なので}} + \ln(g(\lambda_k)) + \frac{1}{s} \lambda^T x + \ln(p(c_k))$$

k に依存しない部分なので、

$$p(c_k|x) = \frac{\exp(\alpha_k)}{\sum_j \exp(\alpha_j)} \quad (4.62)$$

に k 依存しない部分に消滅する。

お金の部分を除いたものを新たに α'_k とおくと、

$$\alpha'_k = \frac{1}{s} \lambda^T x + \ln(g(\lambda_k)) + \ln p(c_k)$$

これは、 x の線形関数なので、事後分布は一般化線形モデルとなる。

(所感)

1. $u(x) = x$ の部分クラスについて論じている点で、指数型分布族一般における議論ではないのではないのか。現に、ガウス分布は $u(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}$ (p.113 (2.221)) であり、共分散行列が共通でない限り線形な境界面とならない。
2. 尺度パラメータの導入は必要あったのか。最後の形を見る限り、 s がなくとも最終形の線形は言えそう

4.3 確率的識別モデル

「4.1 識別関数」では、入力ベクトルから直接識別関数を構築する方法を学んだ。(ref. 最小二乗法、フィッシャーの判別、パーセプトロニアルゴリズム) 「4.2 確率的生成モデル」では、クラスの条件付き確率密度と事前確率を最尤推定し、事後分布を導出するような、生成モデルのアプローチを学んだ。

本節では、 $P(C_k | x)$ を一般化線形モデルで陽に仮定し、最尤推定で直接パラメータを決定する手法を学ぶ。

おさらい：4.1、4.2、4.3の内容はざっくりと以下のイメージ

$$(4.1) \quad y(x) \geq (\text{定数}) \text{ なら } C_1, \\ y(x) < (\text{定数}) \text{ なら } C_2$$

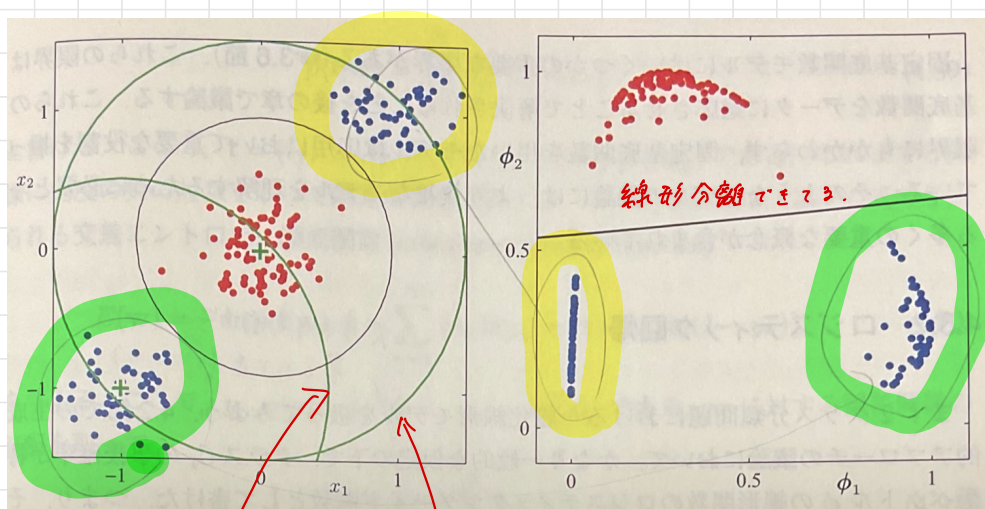
$$(4.2) \quad P(x | C_k) \text{ の分布 を 仮定し, } P(C_k | x) \text{ を 求める,} \\ P(C_k) \\ P(x) \text{ を 導出 でき, データ を 生成 できるようで, 「生成モデル」}$$

$$(4.3) \quad P(C_k | x) \text{ のモデル を 仮定し, 最尤推定 でおこなう.}$$

4.3.1 固定基底関数

本項以降では、入力ベクトル \mathbf{x} の代わりに、基底関数ベクトル $\phi(\mathbf{x})$ を用いることとする。3章の回帰問題では基底関数を導入することで、 \mathbf{x} に対して非線形な目的変数も線形モデルで扱うことができたが、これと同様のアプローチである。

下図は、非線形な基底関数を用いることで線形な識別モデルで分離可能となるような例を示している。



ϕ_1 のガウス基底の等高線。緑の群は ϕ_1 が大きい。黄色の群は ϕ_1 が小さい。

ϕ_2 のガウス基底の等高線。赤の群は ϕ_2 が大きい。

4.3.2 ロジスティック回帰

2クラス分類問題において、4.2節ではいくつかの仮定の元で、事後確率がロジスティックシグモイド関数を用いた一般化線形モデルとなる例を見てきた。

本項では、特に仮定を置くことなく $p(c, \phi)$ をロジスティックシグモイド関数を用いた一般化線形モデルで表現し、最尤推定することを考える。この方法を **ロジスティック回帰** と呼ぶ。

下記のモデル化を行う。

$$p(c, \phi) = y(\phi) = \underbrace{\sigma(w^T \phi)}_{\text{ロジスティックシグモイド関数}}$$

尤度関数は次のようになる。

$$\underbrace{p(\mathbf{t} | \mathbf{w})}_{\text{E 解ラベルベクトル}} = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

これを最大化するのは、負の対数尤度を誤差関数として最小化するのと等価である。この時の誤差関数を **交差エントロピー誤差関数** と呼ぶ。

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

この勾配を考える。4.3.4で述べるが解析的に「勾配=0」の解は得られないので、ここで求める勾配は反復更新に使われる（後述）。

$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} &= -\sum_{n=1}^N \left\{ t_n \frac{1}{y_n} \frac{\partial y_n}{\partial \mathbf{w}} + \frac{1-t_n}{1-y_n} \left(-\frac{\partial y_n}{\partial \mathbf{w}} \right) \right\} \\ &= -\sum_{n=1}^N \left\{ \frac{t_n}{y_n} - \frac{1-t_n}{1-y_n} \right\} \frac{\partial \sigma(\mathbf{w}^T \phi)}{\partial \mathbf{w}} \\ &= +\sum_{n=1}^N \left\{ -\frac{t_n}{y_n} + \frac{1-t_n}{1-y_n} \right\} \underbrace{\sigma(\mathbf{w}^T \phi)(1 - \sigma(\mathbf{w}^T \phi))}_{\text{(演習 4.12) } \frac{d\sigma(a)}{da} = \sigma(a)\{1 - \sigma(a)\} \text{ を使用}} \frac{\partial (\mathbf{w}^T \phi)}{\partial \mathbf{w}} \end{aligned}$$

$$\begin{aligned}
&= -\sum_{n=1}^N \left\{ -\frac{t_n}{y_n} + \frac{1-t_n}{1-y_n} \right\} y_n(1-y_n) \phi \\
&= -\sum_{n=1}^N \left\{ -t_n + t_n y_n + y_n - t_n y_n \right\} \phi \\
&= \sum_{n=1}^N (y_n - t_n) \phi \quad (4.91)
\end{aligned}$$

この勾配を用いて逐次アルゴリズムを構築することで誤差の最小化（つまりここでは最尤推定）を行う。

線形分離可能なデータ集合に対しては、 w の大きさが無限に発散することになる。これを防ぐために w のMAP解を見つければよく、これは誤差関数に正則化項を付加することと等価である。

(ref. p30 (1.67)式)

MAP推定では、 $p(t|w)$ の代わりに、 $p(t|w) p(w)$ と最大化する。 $p(w)$ をガウス分布と仮定し、2次の正則化項が出てきた。p30. 復習

(演習 4.12)

$$\frac{d\sigma(a)}{da} = \sigma(a)(1-\sigma(a)) \quad \text{を示す.}$$

$$\begin{aligned}
\frac{d}{da} \left(\frac{1}{1+e^{-a}} \right) &= \frac{-(-e^{-a})}{(1+e^{-a})^2} \\
&= \frac{1}{1+e^{-a}} \cdot \frac{e^{-a}}{1+e^{-a}} \\
&= \frac{1}{1+e^{-a}} \cdot \frac{1+e^{-a} - 1}{1+e^{-a}} \\
&= \sigma(a) \{1 - \sigma(a)\}
\end{aligned}$$

4.3.3 反復重み付け最小二乗

ロジスティック回帰の誤差関数は、これまでの二次の誤差関数と異なり解析的に最小解が得られない。ただ、誤差関数は凸関数であるため、反復最適化手順によって最小解を得ることができる。

ニュートン法は次の式によってパラメータを更新する。

$$w^{(n(w))} = w^{(old)} - H^{-1} \nabla E(w)$$

線形回帰の最小二乗法の場合、一度の更新で解が得られる。

$$\begin{aligned} E(w) &= \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi_n\}^2 \\ \nabla E(w) &= \sum_{n=1}^N (w^T \phi_n - t_n) \phi_n \\ &= \sum_{n=1}^N \phi_n (\phi_n^T w - t_n) \\ &= \phi_1 \phi_1^T w - \phi_1 t_1 \\ &\quad + \phi_2 \phi_2^T w - \phi_2 t_2 \\ &\quad \vdots \\ &\quad + \phi_n \phi_n^T w - \phi_n t_n \\ &= (\phi_1 \ \phi_2 \ \dots \ \phi_n) \underbrace{\begin{pmatrix} \phi_1^T \\ \vdots \\ \phi_n^T \end{pmatrix}}_{\equiv \Phi} w - (\phi_1 \ \dots \ \phi_n) \underbrace{\begin{pmatrix} t_1 \\ \vdots \\ t_n \end{pmatrix}}_{\equiv t} \\ &= \Phi^T w - \Phi t \\ H = \nabla \nabla E(w) &= \Phi^T \Phi \end{aligned}$$

よって更新式は、

$$\begin{aligned}w^{(new)} &= w^{(old)} - (\Phi^T \Phi)^{-1} \{\Phi^T \Phi w^{(old)} - \Phi^T t\} \\&= w^{(old)} - w^{(old)} + (\Phi^T \Phi)^{-1} \Phi^T t \\&= (\Phi^T \Phi)^{-1} \Phi^T t \quad (\text{正規方程式の解})\end{aligned}$$

ロジスティック回帰の場合は最終的に次のような更新式になる。

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

ここで、 R は $R_{nn} = y_n(1-y_n)$ の対角行列

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

これを以下より示す。

まず、 $\nabla E(w)$ と H を求める。

$$\begin{aligned}\nabla E(w) &= \sum_{n=1}^N (y_n - t_n) \phi_n \quad (\text{前項で求めた}) \\&= \begin{aligned} &\phi_1(y_1 - t_1) \\ &+ \phi_2(y_2 - t_2) \\ &\vdots \\ &+ \phi_n(y_n - t_n) \end{aligned} \\&= \Phi^T (y - t)\end{aligned}$$

$$H = \nabla \nabla E(w) = \frac{\partial}{\partial w} \left\{ \sum_{n=1}^N (y_n - t_n) \phi_n^T \right\}$$
$$\left(\begin{array}{ccc} \frac{\partial^2}{\partial w_1 \partial w_1} & \dots & \frac{\partial^2}{\partial w_1 \partial w_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial w_n \partial w_1} & & \frac{\partial^2}{\partial w_n \partial w_n} \end{array} \right) = \frac{\partial}{\partial w} \left(\frac{\partial}{\partial w_1} \quad \frac{\partial}{\partial w_2} \quad \dots \quad \frac{\partial}{\partial w_n} \right)$$

ここで、 ϕ_n に転置を付けた。

$$\begin{aligned}
&= \sum_{n=1}^N \frac{\partial}{\partial w} y_n \phi_n^T \\
&= \sum_{n=1}^N \frac{\partial}{\partial w} \{ \sigma(w^T \phi) \} \phi_n^T \\
&= \sum_{n=1}^N \sigma(w^T \phi_n) \{ 1 - \sigma(w^T \phi_n) \} \phi_n \phi_n^T \\
&= \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T \\
&= y_1 (1 - y_1) \phi_1 \phi_1^T \\
&\quad + y_2 (1 - y_2) \phi_2 \phi_2^T \\
&\quad \vdots \\
&\quad + y_N (1 - y_N) \phi_N \phi_N^T \\
&= (\underbrace{\phi_1 \dots \phi_N}_{\equiv \Phi^T}) \underbrace{\begin{pmatrix} y_1(1-y_1) & & 0 \\ & \ddots & \\ 0 & & y_N(1-y_N) \end{pmatrix}}_{\equiv R} \underbrace{\begin{pmatrix} \phi_1^T \\ \vdots \\ \phi_N^T \end{pmatrix}}_{\equiv \Phi} \\
&= \Phi^T R \Phi
\end{aligned}$$

ここでのヘッセ行列は、先ほどと違い w に依存する行列となっている。つまり w を更新するたびに H も更新する必要がある。これが、反復「再重み付け」最小二乗と呼ばれる所以である。

また、誤差の2次微分 H が正定値行列であることから、誤差が凸関数であることも示される。(演習4.15)

ここでは、 H が正定値行列であることを示す。
零ベクトルでない任意のベクトル u に対して、

$$\begin{aligned}
u^T H u &= u^T \Phi^T R \Phi u \\
&= u^T R u \\
&= (u_1 \dots u_N) \begin{pmatrix} y_1(1-y_1) & & \\ & \ddots & \\ & & y_N(1-y_N) \end{pmatrix} \begin{pmatrix} u_1 \\ \vdots \\ u_N \end{pmatrix} \\
&= \sum_{n=1}^N u_n^2 y_n (1 - y_n) \\
&> 0 \quad (\because 0 < y_n < 1)
\end{aligned}$$

更新式は次のようになる。

$$\begin{aligned}w^{(new)} &= w^{(old)} - (\Phi^T R \Phi)^{-1} \Phi^T (y - t) \\&= (\Phi^T R \Phi)^{-1} \{ \Phi^T R \Phi w^{(old)} - \Phi^T (y - t) \} \\&= (\Phi^T R \Phi)^{-1} \Phi^T R \underbrace{\{ \Phi w^{(old)} - R^{-1}(y - t) \}}_{= z} \\&= (\Phi^T R \Phi)^{-1} \Phi^T R z\end{aligned}$$

この更新を適用する手法を、反復再重み付け最小二乗法（iterative reweighted least squares method, IRLS）と呼ぶ。

（余談）

ロジスティック回帰をニュートン法で解くという話は、実は「自然科学の統計学」でも少し登場していた（p.238）。ロジスティック回帰はニュートン法などの反復法で解くこと、対数尤度が凹関数であることについて触れている。

演習4.15

$H = \nabla \nabla E(w)$ が正定値行列ならば、 $E(w)$ が凸関数であることを示す。

任意の $w_a, w_b \in \mathbb{R}^m$, $\lambda \in (0, 1)$ に対して.

$$\begin{aligned} f(\lambda) &= E(\lambda w_a + (1-\lambda) w_b) \\ &> \lambda E(w_a) + (1-\lambda) E(w_b) \end{aligned} \quad (*)$$

を言えれば良い。これは f の凸性を示すのと等しい。

$$\frac{df}{d\lambda} = \nabla E(\lambda w_a + (1-\lambda) w_b)^T (w_a - w_b)$$

$$\begin{aligned} \frac{df}{d\lambda^2} &= (w_a - w_b)^T \nabla \nabla E(\lambda w_a + (1-\lambda) w_b) (w_a - w_b) \\ &> 0 \end{aligned}$$

よって f の凸性が言えるので、 $(*)$ を示せた。

4.3.4 多クラスロジスティック回帰

多クラスに拡張した場合も、2クラス同様にモデルを構築することができる。尤度から交差エントロピー誤差関数を定義して、勾配とヘッセ行列から反復更新の手順を得る。

まずは、データが得られた際の各クラスの確率 $p(c_k | \phi)$ を次のようにモデル化する。

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

$$\text{ここで, } a_k = w_k^T \phi$$

負の対数尤度 (交差エントロピー誤差関数) は次のようになる。

$$\begin{aligned} E(w_1, \dots, w_K) &= -\ln p(\pi | w_1, \dots, w_K) \\ &= -\ln \prod_{n=1}^N \prod_{k=1}^K p(c_k | \phi_n)^{t_{nk}} \\ &= -\ln \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \\ &= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \end{aligned}$$

この w_i に関する勾配は次のようになる。

$$\begin{aligned} \frac{\partial}{\partial w_i} E(w_1, \dots, w_K) &= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \frac{1}{y_{nk}} \frac{\partial y_{nk}}{\partial a_i} \frac{\partial a_i}{\partial w_i} \\ &= -\sum_{n=1}^N \sum_{k=1}^K \frac{t_{nk}}{y_{nk}} y_{nk} (I_{ki} - y_{ni}) \phi_n \\ &\quad \text{(演習 4.17)} \\ &= -\sum_{n=1}^N \sum_{k=1}^K (t_{nk} I_{ki} - t_{nk} y_{ni}) \phi_n \end{aligned}$$

$$= - \sum_{n=1}^N \left(\sum_{k=1}^K t_{nk} I_{kj} - \sum_{k=1}^K t_{nk} y_{nj} \right) \phi_n$$

$\underbrace{\hspace{10em}}_{k=j \text{ の } 2 \text{ 項}} - \underbrace{\hspace{10em}}_{=1}$
 残差

$$= - \sum_{n=1}^N (t_{nj} - y_{nj}) \phi_n$$

$$= \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

(演習4.17)

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

$$\frac{\partial y_k}{\partial a_j} = \frac{\partial}{\partial a_j} \frac{\exp(a_k)}{\sum_i \exp(a_i)}$$

$$\frac{\partial}{\partial \alpha_j} \exp(\alpha_k) = \frac{\partial \alpha_k}{\partial \alpha_j} \cdot \frac{\partial}{\partial \alpha_k} \exp(\alpha_k) = I_{kj} \exp(\alpha_k)$$

ヘッセ行列について考える前に、多クラスロジスティック回帰におけるニュートン法の更新式を考える。PRMLには詳しく載っていないので、こちらを参考にした：

<https://www.iwanttobeacat.com/entry/2018/09/08/010518>

結論としては、次の更新式を使う。

$$\begin{pmatrix} (KM, 1) \\ w_1^{(new)} \\ w_2^{(new)} \\ \vdots \\ w_K^{(new)} \end{pmatrix} = \begin{pmatrix} (KM, 1) \\ w_1^{(old)} \\ w_2^{(old)} \\ \vdots \\ w_K^{(old)} \end{pmatrix} - \begin{pmatrix} (KM, KM) \\ H_{11} & \dots & H_{1K} \\ \vdots & & \vdots \\ H_{K1} & \dots & H_{KK} \end{pmatrix}^{-1} \begin{pmatrix} (KM, 1) \\ \frac{\partial E}{\partial w_1^{(old)}} \\ \vdots \\ \frac{\partial E}{\partial w_K^{(old)}} \end{pmatrix}$$

↑
1個1個が $M \times M$ 行列
 M は w の次元 (p204)

ここで、ヘッセ行列のうち、ブロック j, k について次のように求めることができる。

$$\begin{aligned} \nabla_{w_j} \nabla_{w_k} E(w_1, \dots, w_K) &= \frac{\partial}{\partial w_j} \frac{\partial}{\partial w_k^T} E \\ &= \frac{\partial}{\partial w_j} \sum_{n=1}^N (y_{nk} - t_{nk}) \phi_n^T \\ &= \sum_{n=1}^N \frac{\partial y_{nk}}{\partial a_j} \frac{\partial a_j}{\partial w_j} \phi_n^T \\ &= \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \phi_n \phi_n^T \\ &= (\phi_1, \dots, \phi_n) \begin{pmatrix} y_{1k}(I_{kj} - y_{1j}) & & 0 \\ & \ddots & \\ 0 & & y_{nk}(I_{kj} - y_{nj}) \end{pmatrix} \begin{pmatrix} \phi_1^T \\ \vdots \\ \phi_n^T \end{pmatrix} \\ &= \Phi^T R_{k,j} \Phi \end{aligned}$$

教科書と演算子の
順序違うけど。
こいつよね？