

パターン認識と機械学習

第4章 線形識別モデル

橋本龍二

January 28, 2021

分類問題に関する3つのアプローチ：

1. 識別関数

入力ベクトル \mathbf{x} から直接クラスを推定する関する $f(\mathbf{x})$ を見つける。2 クラスなら $f(\cdot)$ は2値を出力する。

2. 推論段階で条件付き確率 $p(\mathcal{C}_k|\mathbf{x})$ をモデル化したのち、 $p(\mathcal{C}_k|\mathbf{x})$ を利用して決定理論を適用する。決定理論は、例えば

$$\min_j \sum_k L_{kj} p(\mathcal{C}_k|\mathbf{x})$$

となるように j を決める期待損失最小化や、

$$\begin{aligned} \max \sum_k p(\mathbf{x} \in \mathcal{R}_k, C_k) &= \max \sum_k \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x} \\ &= \max \sum_k \int_{\mathcal{R}_k} p(\mathbf{x}) p(\mathcal{C}_k|\mathbf{x}) d\mathbf{x} \end{aligned}$$

となるように決定領域 \mathcal{R}_k を決める誤識別率最小化がある。
推論段階において、 $p(\mathcal{C}_k|\mathbf{x})$ を決める方法には次の二通りがある。

(a) パラメトリックに分布を表現する。(確率的識別モデル)

(b) クラスの事前確率 $p(\mathcal{C}_k)$ とクラスで条件づけられた確率密度 $p(\mathbf{x}|\mathcal{C}_k)$ から各クラスの事後確率：

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathcal{C}_k)p(\mathbf{x}|\mathcal{C}_k)}{p(\mathbf{x})}$$

を求める。(確率的生成モデル)

以上の1、2(a)、2(b)の方法について詳しく見ていく。

1 識別関数

1.1 2 クラス

D 次元の入力ベクトル \mathbf{x} が与えられた時、線形識別関数

$$y(\mathbf{x}) = {}^T \mathbf{w} \mathbf{x} + w_0$$

が 0 以上であればクラス \mathcal{C}_1 に割り当てるこにする (決定境界 : $y(\mathbf{x}) = 0$ を満たす点 \mathbf{x} の集合)。

$${}^T \mathbf{w} \mathbf{x} + w_0 = 0$$

を満たす決定境界は $D - 1$ 次元の超平面である ($D - 1$ 個の値が決まれば残りは方程式より決定する)。

決定境界上の任意の 2 点 $\mathbf{x}_A, \mathbf{x}_B$ を考える。

$${}^T \mathbf{w} (\mathbf{x}_A - \mathbf{x}_B) = 0 \quad (1)$$

より、ベクトル \mathbf{w} は決定面上の任意のベクトルと直交する。 $(\mathbf{w}$ は平面の法線ベクトル、みたいなこと)

今、決定面上の任意の点 \mathbf{x} をとり、原点から決定面への垂線をひく。この垂線は \mathbf{w} と平行なので、定数 t を用いて $t\mathbf{w}$ と表すことにする。 $t\mathbf{w}$ の長さが原点から決定面への距離である。(1) より

$$\begin{aligned} (\mathbf{x} - t\mathbf{w}, \mathbf{w}) &= 0 \\ (\mathbf{x}, \mathbf{w}) - t\|\mathbf{w}\|^2 &= 0 \\ \therefore t &= \frac{(\mathbf{x}, \mathbf{w})}{\|\mathbf{w}\|^2} \end{aligned}$$

したがって、原点から決定面の距離は、

$$\begin{aligned} \|t\mathbf{w}\| &= \left\| \frac{(\mathbf{x}, \mathbf{w})}{\|\mathbf{w}\|^2} \mathbf{w} \right\| = \frac{(\mathbf{x}, \mathbf{w})}{\|\mathbf{w}\|} \\ &= -\frac{w_0}{\|\mathbf{w}\|} \end{aligned}$$

(ただし、 \mathbf{w} がどちら向きでも r は変わらない訳だから厳密に距離を求めたくなったら絶対値を取る必要がある。)

次に、任意の点 \mathbf{x} と決定面との距離を求める。 \mathbf{x} から決定面へ下ろした垂線と決定面との交点を \mathbf{x}_\perp (\mathbf{x} の決定面への直交射影) と置くと、垂線は \mathbf{w} と平行なので、定数 s を用いて $s\mathbf{w}$ と表すことにする。すると \mathbf{x} は、

$$\begin{aligned} \mathbf{x} &= \mathbf{x}_\perp + s\mathbf{w} \\ &= \mathbf{x}_\perp + \frac{r}{\|\mathbf{w}\|} \mathbf{w} \quad (2) \\ &\quad (r = \|\mathbf{w}\|s) \end{aligned}$$

と書ける。 r が \mathbf{x} と決定面の距離である。今、(2) の両辺に ${}^T \mathbf{w}$ を掛け w_0 を加えることで、

$$\begin{aligned} {}^T \mathbf{w} \mathbf{x} + w_0 &= {}^T \mathbf{w} (\mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}) + w_0 \\ &= {}^T \mathbf{w} \mathbf{x}_\perp + w_0 + r \|\mathbf{w}\| \\ \therefore y(\mathbf{x}) &= r \|\mathbf{w}\| \\ \therefore r &= \frac{y(\mathbf{x})}{\|\mathbf{w}\|} \end{aligned}$$

まとめ

$$y(\mathbf{x}) = {}^T \mathbf{w} \mathbf{x} + w_0$$

なる線形識別関数において、

- \mathbf{w} は決定面と直交する。すなわち \mathbf{w} は決定面の方向を決定する。
- 原点から決定面までの距離は $-\frac{w_0}{\|\mathbf{w}\|}$ である。すなわち w_0 は決定面の位置を決定する。
- 任意の点 \mathbf{x} と決定面との距離は $\frac{y(\mathbf{x})}{\|\mathbf{w}\|}$ である。すなわち $y(\mathbf{x})$ は決定面からの点 \mathbf{x} の距離を与える。

1.2 多クラス

D 次元入力ベクトル \mathbf{x} に対し、

$$y_k(\mathbf{x}) = {}^T \mathbf{w}_k \mathbf{x} + w_{k0}$$

を定義し、全ての $j \neq k$ に対し $y_k(\mathbf{x}) > y_j(\mathbf{x})$ ならば \mathbf{x} をクラス \mathcal{C}_k に分類することを考える。この時、クラス \mathcal{C}_k とクラス \mathcal{C}_j の決定境界は $y_k(\mathbf{x}) = y_j(\mathbf{x})$ である。すなわち、

$${}^T (\mathbf{w}_k - \mathbf{w}_j) \mathbf{x} = -(w_{k0} - w_{j0})$$

なる $D - 1$ 次元超平面。

識別器 $y_k(\mathbf{x})$ の決定領域は必ず凸領域である。つまり、任意の $\mathbf{x}_A, \mathbf{x}_B \in \mathcal{R}_k$ を取れば、線分 $\mathbf{x}_A \mathbf{x}_B$ 上の任意の点 $\hat{\mathbf{x}}$ は必ず \mathcal{R}_k に属する。なぜなら、

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B, \quad 0 \leq \lambda \leq 1$$

とすれば、

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda)y_k(\mathbf{x}_B)$$

今、任意の $j \neq k$ について $y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A), y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B)$ が成り立つのだから、 $\hat{\mathbf{x}} \in \mathcal{R}_k$

2 クラスの場合でも、 $y_1(\mathbf{x}), y_2(\mathbf{x})$ の 2 つの識別関数によって定式化することも可能である。

1.3 分類における最小二乗

D 次元入力 \mathbf{x} が与えられた時、2 値表記法の目標ベクトル \mathbf{t} を $\mathbf{y}(\mathbf{x})$ で推定する。最小二乗法で $\mathbf{y}(\mathbf{x}) = E[\mathbf{t}|\mathbf{x}]$ とすれば、 \mathbf{y} の第 k 要素は、クラス C_k に割り当てる確率である（識別関数なので厳密には確率ではないことに注意。確率っぽく扱う）。今、 n 個目のデータ \mathbf{x}_n を C_k に割り当てる確率を

$$y_k(\mathbf{x}) = {}^T \mathbf{w}_k \mathbf{x} + w_{k0}, \quad k = 1, \dots, K$$

とすれば、 \mathbf{x}_n に対する推定値は

$$\begin{aligned} \mathbf{y}(\mathbf{x}_n) &= \begin{pmatrix} y_1(\mathbf{x}_n) \\ \vdots \\ y_K(\mathbf{x}_n) \end{pmatrix} \\ &= \begin{pmatrix} {}^T \tilde{\mathbf{w}}_1 \tilde{\mathbf{x}}_n \\ \vdots \\ {}^T \tilde{\mathbf{w}}_K \tilde{\mathbf{x}}_n \end{pmatrix} \\ &= {}^T \tilde{W} \tilde{\mathbf{x}}_n \end{aligned}$$

ただし、

$$\tilde{\mathbf{w}}_k = \begin{pmatrix} w_0 \\ \mathbf{w}_k \end{pmatrix}, \quad \tilde{\mathbf{x}}_n = \begin{pmatrix} 1 \\ \mathbf{x}_n \end{pmatrix}, \quad \tilde{W} = \begin{pmatrix} \tilde{\mathbf{w}}_1 \\ \vdots \\ \tilde{\mathbf{w}}_K \end{pmatrix}$$

とできる。 $\mathbf{y}(\mathbf{x}_n), n = 1, \dots, N$ を N 個縦に並べた行列を Y とすれば、

$$\begin{aligned} Y &= \begin{pmatrix} {}^T \mathbf{y}(\mathbf{x}_1) \\ \vdots \\ {}^T \mathbf{y}(\mathbf{x}_N) \end{pmatrix} \\ &= \begin{pmatrix} y_1(\mathbf{x}_1) & \dots & y_K(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ y_1(\mathbf{x}_N) & \dots & y_K(\mathbf{x}_N) \end{pmatrix} \end{aligned}$$

また、 Y は

$$\begin{aligned} Y &= \begin{pmatrix} {}^T\mathbf{y}(\mathbf{x}_1) \\ \vdots \\ {}^T\mathbf{y}(\mathbf{x}_N) \end{pmatrix} \\ &= \begin{pmatrix} {}^T\tilde{\mathbf{x}}_1\tilde{W} \\ \vdots \\ {}^T\tilde{\mathbf{x}}_N\tilde{W} \end{pmatrix} \\ &= \tilde{X}\tilde{W}, \quad (\tilde{X} = \begin{pmatrix} {}^T\tilde{\mathbf{x}}_1 \\ \vdots \\ {}^T\tilde{\mathbf{x}}_N \end{pmatrix}) \end{aligned}$$

である。二乗和誤差関数を求める。(二重和はトレースで書ける！)

$$\begin{aligned} E_D(\tilde{W}) &= \frac{1}{2} \sum_{k=1}^K \sum_{n=1}^N \{y_k(\mathbf{x}_n) - t_{nk}\}^2 \\ &= \frac{1}{2} Tr \left(\begin{matrix} (y_1(\mathbf{x}_1) - t_{11})^2 + \dots + (y_1(\mathbf{x}_N) - t_{N1})^2 \\ \ddots \\ (y_K(\mathbf{x}_1) - t_{1K})^2 + \dots + (y_K(\mathbf{x}_N) - t_{NK})^2 \end{matrix} \right) \\ &= \frac{1}{2} Tr \left(\begin{matrix} y_1(\mathbf{x}_1) - t_{11} & \dots & y_1(\mathbf{x}_N) - t_{N1} \\ \vdots & \ddots & \vdots \\ y_K(\mathbf{x}_1) - t_{1K} & \dots & y_K(\mathbf{x}_N) - t_{NK} \end{matrix} \right) \left(\begin{matrix} y_1(\mathbf{x}_1) - t_{11} & \dots & y_K(\mathbf{x}_1) - t_{1K} \\ \vdots & \ddots & \vdots \\ y_1(\mathbf{x}_N) - t_{N1} & \dots & y_K(\mathbf{x}_N) - t_{NK} \end{matrix} \right) \\ &= \frac{1}{2} \{ {}^T(Y - T)(Y - T) \} \\ &= \frac{1}{2} Tr \{ {}^T(\tilde{X}\tilde{W} - T)(\tilde{X}\tilde{W} - T) \} \end{aligned}$$

二乗和誤差関数を \tilde{W} で微分して 0 と置くと、

$$\tilde{W} = ({}^T\tilde{X}\tilde{X})^{-1} {}^T\tilde{X}T = \tilde{X}^\dagger T$$

なので、求める識別関数は

$$\mathbf{y}(\mathbf{x}) = {}^T\tilde{W}\tilde{x} = {}^T T {}^T\tilde{X}^\dagger \tilde{x}$$

この関数は任意の \mathbf{x} に対し各要素の和が 1 になる、という性質を持つ（したがって確率かのように取り扱えるが、値域は $(0, 1)$ でない）。

最小二乗法の弱点（回帰でも同じことが言える）：

- 最小二乗法は条件付きガウス分布のもとでの最尤推定に相当するので、外れ値に頑健でない。
- 最小二乗法は条件付き分布 $p(t|x)$ が正規分布しない場合仮定を満たさないのでからうまくいかない。

1.4 フィッシャーの線形判別

D 次元入力ベクトルを

$$y = {}^T \mathbf{w} \mathbf{x}$$

によって 1 次元に射影して $y > -w_0$ の時クラス \mathcal{C}_1 に分類する二分類問題を考える。クラス \mathcal{C}_1 の点が N_1 個、クラス \mathcal{C}_2 の点が N_2 個あるとすると、各クラスの平均ベクトルは

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n$$

である。射影先で分離度が最大になるように \mathbf{w} を定めたい。

$$\max_{\mathbf{w}} \quad m_2 - m_1 = {}^T \mathbf{w} (\mathbf{m}_2 - \mathbf{m}_1) \quad s.t. \quad \sum_i w_i^2 = 1$$

一階条件より、

$$\mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1 \quad \left(\sum_i w_i^2 = 1 \right)$$

より、 \mathbf{x} 一次元数ベクトル空間 $\langle \mathbf{w} \rangle = \langle \mathbf{m}_2 - \mathbf{m}_1 \rangle$ へと射影 $(\mathbf{w}, \mathbf{x})\mathbf{w}$ してそこで閾値を決めれば良いのだが、図 4.6 左のようにこれではまだ射影先でクラスの被りが見られる。これを防ぐために、射影先での各クラスでのクラス内分散を

$$s_k = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2$$

と定義して、

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

を最大化する問題を考える。今、

$$\begin{aligned}
 (m_2 - m_1)^2 &= ({}^T \mathbf{w} \mathbf{m}_2 - {}^T \mathbf{w} \mathbf{m}_1)^2 \\
 &= {}^T \mathbf{w} (\mathbf{m}_2 - \mathbf{m}_1) {}^T (\mathbf{m}_2 - \mathbf{m}_1) \mathbf{w} \\
 s_1^2 + s_2^2 &= \sum_{n \in \mathcal{C}_1} ({}^T \mathbf{w} \mathbf{x}_n - {}^T \mathbf{w} \mathbf{m}_1)^2 + \sum_{n \in \mathcal{C}_2} ({}^T \mathbf{w} \mathbf{x}_n - {}^T \mathbf{w} \mathbf{m}_2)^2 \\
 &= {}^T \mathbf{w} \left\{ \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1) {}^T (\mathbf{x}_n - \mathbf{m}_1) + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2) {}^T (\mathbf{x}_n - \mathbf{m}_2) \right\} \mathbf{w}
 \end{aligned}$$

なので、

$$\begin{aligned}
 S_B &= (\mathbf{m}_2 - \mathbf{m}_1) {}^T (\mathbf{m}_2 - \mathbf{m}_1) \\
 S_W &= \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1) {}^T (\mathbf{x}_n - \mathbf{m}_1) + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2) {}^T (\mathbf{x}_n - \mathbf{m}_2)
 \end{aligned}$$

とおけば、最大化問題は

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{{}^T \mathbf{w} S_B \mathbf{w}}{{}^T \mathbf{w} S_W \mathbf{w}} \quad s.t. \quad \sum_i w_i^2 = 1$$

と書ける。一階条件を解くことで、

$$\begin{aligned}
 \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \frac{{}^T \mathbf{w} S_B \mathbf{w}}{{}^T \mathbf{w} S_W \mathbf{w}} \\
 &= \frac{2({}^T \mathbf{w} S_W \mathbf{w}) S_B \mathbf{w} - 2({}^T \mathbf{w} S_B \mathbf{w}) S_W \mathbf{w}}{({}^T \mathbf{w} S_W \mathbf{w})^2} = \mathbf{0} \\
 \therefore ({}^T \mathbf{w} S_B \mathbf{w}) S_W \mathbf{w} &= ({}^T \mathbf{w} S_W \mathbf{w}) S_B \mathbf{w} \tag{1}
 \end{aligned}$$

今、

$$S_B \mathbf{w} = (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1, \mathbf{w})$$

より、 $S_B \mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1$ である。また(1)においてスカラー部分 ${}^T \mathbf{w} S_B \mathbf{w}$, ${}^T \mathbf{w} S_W \mathbf{w}$ を無視すれば、

$$\begin{aligned}
 \mathbf{w} &\propto S_W^{-1} S_B \mathbf{w} \\
 &\propto S_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)
 \end{aligned}$$

クラス内共分散が等方的（きれいな円形であればどの断面で切っても分散は変わらない）であれば S_W は単位行列に比例し、 $\mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1$ となる。

2 確率的生成モデル

二クラス分類においてデータ \mathbf{x} が与えられたときクラス \mathcal{C}_1 の事後確率を計算する。

$$\begin{aligned}
 p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathcal{C}_1 \wedge \mathbf{x})}{p(\mathbf{x})} \\
 &= \frac{p(\mathcal{C}_1 \wedge \mathbf{x})}{p(\mathcal{C}_1 \wedge \mathbf{x}) + p(\mathcal{C}_2 \wedge \mathbf{x})} \\
 &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\
 &= \frac{1}{1 + \frac{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}} \tag{1}
 \end{aligned}$$

ここで

$$\begin{aligned}
 \frac{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)} &= \exp\{\log p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2) - \log p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)\} \\
 &= \exp\{-\log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}\}
 \end{aligned}$$

なので、

$$\begin{aligned}
 (1) &= \frac{1}{1 + \exp(-a)}, \quad a = \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\
 &= \sigma(a)
 \end{aligned}$$

このように、二クラス分類の事後確率はロジスティックシグモイド関数 σ で書ける。 σ の逆関数：

$$a = \log\left(\frac{\sigma}{1-\sigma}\right) = \log \frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})}$$

をロジット関数という。 $p(\mathcal{C}_1|\mathbf{x}) \geq p(\mathcal{C}_2|\mathbf{x})$ の時 $a > 0$ となり、 $0.5 \leq \sigma \leq 1$ となる。

2.1 連続値入力

クラス \mathcal{C}_k の確率密度を

$$\begin{aligned}
 p(\mathbf{x}|\mathcal{C}_k) &= N(\boldsymbol{\mu}_k, \Sigma) \\
 &= \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} {}^T(\mathbf{x} - \boldsymbol{\mu}_k) \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}
 \end{aligned}$$

と定義する。2 クラス分類においてロジット関数は

$$\begin{aligned}
a &= \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\
&= -\frac{1}{2}\left\{ {}^T(\mathbf{x} - \boldsymbol{\mu}_1)\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) - {}^T(\mathbf{x} - \boldsymbol{\mu}_1)\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_2)\right\} + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)} \\
&= -\frac{1}{2}\left\{ {}^T\boldsymbol{\mu}_1\Sigma^{-1}\boldsymbol{\mu}_1 - {}^T\boldsymbol{\mu}_2\Sigma^{-1}\boldsymbol{\mu}_2 - 2 {}^T\mathbf{x}\Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\right\} + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)} \\
&= {}^T\mathbf{x}\Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2} {}^T\boldsymbol{\mu}_1\Sigma^{-1}\boldsymbol{\mu}_1 + \frac{1}{2} {}^T\boldsymbol{\mu}_2\Sigma^{-1}\boldsymbol{\mu}_2 + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}
\end{aligned}$$

となるので、クラス \mathcal{C}_1 に対する事後確率は

$$\begin{aligned}
p(\mathcal{C}_1|\mathbf{x}) &= \sigma(a) \\
&= \sigma({}^T\mathbf{w}\mathbf{x} + w_0)
\end{aligned}$$

ただし、

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad w_0 = -\frac{1}{2} {}^T\boldsymbol{\mu}_1\Sigma^{-1}\boldsymbol{\mu}_1 + \frac{1}{2} {}^T\boldsymbol{\mu}_2\Sigma^{-1}\boldsymbol{\mu}_2 + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$

と置いた。このように、各クラスの事後確率は \mathbf{x} の線形関数のロジスティックシグモイド関数によって与えられる。

2.2 最尤解

2 クラス分類において、 D 次元ベクトル \mathbf{x} とクラスラベル t (クラス \mathcal{C}_1 の時 $t = 1$ 、 \mathcal{C}_2 の時 $t = 0$) として訓練データ $\{\mathbf{x}_n, t_n\}$ が与えられたとする。クラス事前確率を

$$p(\mathcal{C}_1) = \pi, \quad p(\mathcal{C}_2) = 1 - \pi \quad (0 \leq \pi \leq 1)$$

と定めれば

$$\begin{aligned}
p(\mathbf{x}_n, \mathcal{C}_1) &= p(\mathcal{C}_1)P(\mathbf{x}_n|\mathcal{C}_1) = \pi N(\mathbf{x}_n|\boldsymbol{\mu}_1, \Sigma) \\
p(\mathbf{x}_n, \mathcal{C}_2) &= p(\mathcal{C}_2)P(\mathbf{x}_n|\mathcal{C}_2) = (1 - \pi)N(\mathbf{x}_n|\boldsymbol{\mu}_2, \Sigma)
\end{aligned}$$

より、対数尤度関数は

$$\begin{aligned}
p(\mathbf{t}, X|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \prod_{n=1}^N [\pi N(\mathbf{x}_n|\boldsymbol{\mu}_1, \Sigma)]^{t_n} [(1 - \pi)N(\mathbf{x}_n|\boldsymbol{\mu}_2, \Sigma)]^{1-t_n} \\
\therefore \log p(\mathbf{t}, X|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \sum_{n=1}^N \{t_n \log \pi + (1 - t_n) \log (1 - \pi) + \log N(\mathbf{x}_n|\boldsymbol{\mu}_1, \Sigma) \\
&\quad + (1 - t_n) \log N(\mathbf{x}_n|\boldsymbol{\mu}_2, \Sigma)\}
\end{aligned}$$

これを π について微分して 0 と置くことで、

$$\begin{aligned}\frac{\partial}{\partial \pi} \log p(\mathbf{t}, X | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma) &= \sum_{n=1}^N \left(\frac{t_n}{\pi} - \frac{1-t_n}{1-\pi} \right) \\&= \sum_{n=1}^N \frac{t_n - \pi}{\pi(1-\pi)} = 0 \\ \therefore \hat{\pi}^{ML} &= \frac{1}{N} \sum_{n=1}^N t_n \\&= \frac{N_1}{N} \\&= \frac{N_1}{N_1 + N_2}\end{aligned}$$

これでクラス事前確率の最尤推定量が得られた。つぎに μ_1, μ_2 について微分して 0 と置くことで

$$\begin{aligned}
\frac{\partial}{\partial \mu_1} &= \frac{\partial}{\partial \mu_1} \sum_{n=1}^N t_n \log N(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma) \\
&= \frac{\partial}{\partial \mu_1} \left\{ -\frac{1}{2} \sum_{n=1}^N t_n {}^T (\mathbf{x}_n - \boldsymbol{\mu}_1) \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \right\} \\
&= -\frac{1}{2} \sum_{n=1}^N t_n \left(2 \sum_{n=1}^{-1} \boldsymbol{\mu}_1 - 2 \Sigma^{-1} \mathbf{x}_n \right) \\
&= -N_1 \sum_{n=1}^{-1} \boldsymbol{\mu}_1 + \Sigma^{-1} \sum_{n=1}^N t_n \mathbf{x}_n = 0 \\
\therefore \hat{\boldsymbol{\mu}}_1^{ML} &= \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n \\
\frac{\partial}{\partial \mu_2} &= \frac{\partial}{\partial \mu_2} \sum_{n=1}^N (1 - t_n) \log N(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma) \\
&= \frac{\partial}{\partial \mu_2} \left\{ -\frac{1}{2} \sum_{n=1}^N (1 - t_n) {}^T (\mathbf{x}_n - \boldsymbol{\mu}_2) \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \right\} \\
&= -\frac{1}{2} \sum_{n=1}^N (1 - t_n) (2 \Sigma^{-1} \boldsymbol{\mu}_2 - 2 \Sigma^{-1} \mathbf{x}_n) \\
&= -N_2 \Sigma^{-1} \boldsymbol{\mu}_2 + \Sigma^{-1} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n = 0 \\
\therefore \hat{\boldsymbol{\mu}}_2^{ML} &= \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n
\end{aligned}$$

より、 $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2$ は各クラスに割り当てる入力ベクトルの平均。

まとめ：

各クラスについて入力データにガウス分布を仮定すれば、クラス事後確率はガウス分布のパラメータを用いてロジスティックシグモイド関数によって与えられる。そこで、パラメータの最尤推定量を求め、入力に対し各クラスの事後確率を推定することができるので、決定理論を適用してクラス分類を行える。また、 $p(\mathbf{x}|\mathcal{C}_k)$ を用いて各クラスの入力ベクトルをサンプリングすること

も可能である。一方、この方法は外れ値に脆弱である。

3 確率的識別モデル

$p(\mathbf{x}|\mathcal{C}_k)$, $p(\mathcal{C}_k)$ は考えずに直接 $p(\mathcal{C}_k|\mathbf{x})$ をモデル化する。

3.1 固定基底関数

基底関数ベクトル ϕ によって入力ベクトル \mathbf{x} を特徴空間に写像し、特徴空間内で線形にクラス分類することを考える。

3.2 ロジスティック回帰

特徴ベクトル ϕ が得られた時、クラス \mathcal{C}_1 の事後確率は

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(^T w \phi)$$

の形で書ける。ただし、 $p(\mathcal{C}_2|\phi) = 1 - p(\mathcal{C}_1|\phi)$ である。これにより、クラス t は確率 $p(\mathcal{C}_t|\phi)$ のベルヌーイ分布に従う。 $t \sim Ber(p(\mathcal{C}_t|\phi))$ 生成モデルでは、 $p(\phi|\mathcal{C}_k) = N(\mu_k, \Sigma)$ と置いて最尤推定によって w を決める。しかし、特徴空間が M 次元とすると、この場合 μ_1 と μ_2 で $2M$ 個、 Σ には

$$M + (M - 1) + \dots + 1 = \frac{M(M + 1)}{2}$$

で、合わせて $M(M + 5)/2 + 1$ 個ものパラメータが必要になる。対してロジスティック回帰でパラメータを直接推定する場合は調整可能パラメータの数は M 個に抑えられる。

データ集合 $\{\phi_n, t_n\}$ が与えられた時、データ N 個の対数尤度関数は

$$\begin{aligned} p(\mathbf{t}|w) &= \prod_{n=1}^N Ber(t_n|p(\mathcal{C}_1|\phi_n)) \\ &= \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n} \\ \log p(\mathbf{t}|w) &= \sum_{n=1}^N \{t_n \log y_n + (1 - t_n) \log(1 - y_n)\} \end{aligned}$$

ただし、

$$y_n = p(\mathcal{C}_1|\phi_n)$$

と置いた。交差エントロピー誤差とは、対数尤度関数の負で、

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

のことである。 \mathbf{w} に対する誤差関数の勾配をとって

$$\begin{aligned}\nabla E(\mathbf{w}) &= -\sum_{n=1}^N \left\{ \frac{t_n}{y_n} \frac{\partial}{\partial \mathbf{w}} y_n - \frac{1-t_n}{1-y_n} \frac{\partial}{\partial \mathbf{w}} y_n \right\} \\ &= -\sum_{n=1}^N \frac{\partial y_n}{\partial \mathbf{w}} \frac{t_n(1-y_n) - y_n(1-t_n)}{y_n(1-y_n)} \\ &= \sum_{n=1}^N y_n(1-y_n) \frac{\partial^T \mathbf{w} \phi_n}{\partial \mathbf{w}} \frac{y_n - t_n}{y_n(1-y_n)} \\ &= \sum_{n=1}^N (y_n - t_n) \phi_n\end{aligned}$$

逐次学習ならば、 $-\nabla E_n(\mathbf{w})$ の方向に重みベクトルを更新していくべき。しかし、ロジスティックシグモイド関数 y_n の非線形性により、最尤解を解析的に導出することはできない。

3.3 反復再重み付け最小二乗

誤差関数 $E(\mathbf{w})$ を最小化する方法としてニュートン・ラフソン法がある。ニュートン・ラフソン法の \mathbf{w} 更新は、

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

によって行う。ただし、 H は \mathbf{w} に関する $E(\mathbf{w})$ の二階微分を要素とするヘッセ行列。

例：

$$t = {}^T \mathbf{w} \phi(\mathbf{x}) + \epsilon, \quad \epsilon \sim N(0, \beta^{-1})$$

なる線形回帰モデルを考える。ここで、二乗和誤差関数は

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