

はじめてのパターン認識

第6章 線形識別関数

山木翔馬

2015/06/16

2クラス問題の線形識別関数

線形識別関数の定義

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

$$\begin{array}{ll} \mathbf{x} = (x_1, \dots, x_d)^T & d \text{次元入力ベクトル} \\ \mathbf{w} = (w_1, \dots, w_d)^T & \text{係数ベクトル} \\ w_0 & \text{バイアス項} \end{array}$$

識別境界を $f(\mathbf{x}) = 0$ とすれば

識別規則

$$f(\mathbf{x}) = \begin{cases} C_1 & (f(\mathbf{x}) \geq 0) \\ C_2 & (f(\mathbf{x}) < 0) \end{cases}$$

超平面の方程式 (1/3)

識別境界では $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ より

$$f(\mathbf{x}) = 0 \Leftrightarrow \mathbf{w}^T \mathbf{x} = -w_0$$

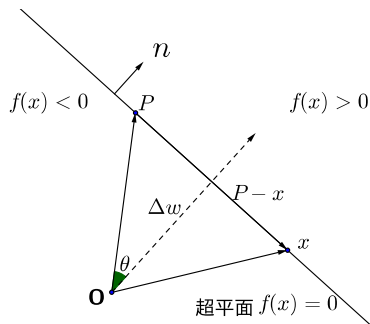
両辺を $\|\mathbf{w}\|$ で正規化, $\mathbf{n} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$, $\Delta_w = -\frac{w_0}{\|\mathbf{w}\|}$

$$\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|} \Leftrightarrow \mathbf{n}^T \mathbf{x} = \Delta_w$$

識別境界

$$f(\mathbf{x}) = \mathbf{n}^T \mathbf{x} - \Delta_w = 0$$

超平面の方程式 (2/3)



識別境界上の任意の点の位置ベクトル \mathbf{P} について

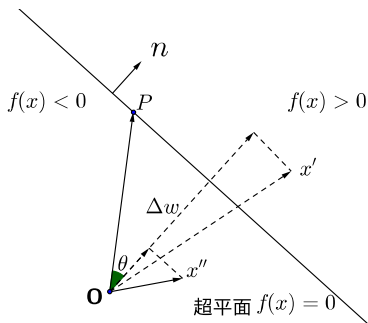
$$f(\mathbf{P}) = \mathbf{n}^T \mathbf{x} - \Delta_w = 0$$

が成り立つので

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{n}^T \mathbf{x} - \Delta_w \\ &= \mathbf{n}^T (\mathbf{x} - \mathbf{P}) = 0 \end{aligned}$$

- $\mathbf{x} - \mathbf{P} \perp \mathbf{n}$
- 識別境界は単位法線ベクトル \mathbf{n} をもつ超平面
- $\Delta_w = \mathbf{n}^T \mathbf{P}$ は原点から識別超平面までの距離
($\Delta_w = \mathbf{n}^T \mathbf{P} = \|\mathbf{n}\| \cdot \|\mathbf{P}\| \cos\theta = \|\mathbf{P}\| \cos\theta$)
- Δ_w : 正規化されたバイアス

超平面の方程式 (3/3)



$$\begin{aligned} & \mathbf{n}^T \mathbf{x}' > \Delta_w \\ \Leftrightarrow & f(\mathbf{x}') = \mathbf{n}^T \mathbf{x}' - \Delta_w > 0 \\ \Rightarrow & \mathbf{x}' \in C_1 \end{aligned}$$

$$\begin{aligned} & \mathbf{n}^T \mathbf{x}'' < \Delta_w \\ \Leftrightarrow & f(\mathbf{x}'') = \mathbf{n}^T \mathbf{x}'' - \Delta_w < 0 \\ \Rightarrow & \mathbf{x}'' \in C_2 \end{aligned}$$

多クラス問題への拡張

- 一対他 (one-versus-the-rest)
- 一対一 (one-versus-one)

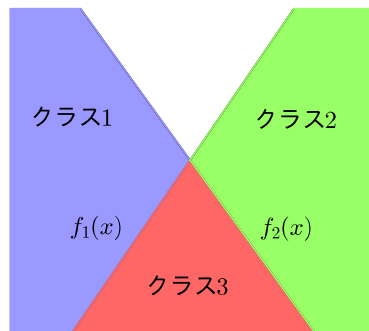
上2つの方法では、クラス決定ができない領域がある。
それを解消したのが

- 最大識別関数法

一対他 (one-versus-the-rest)

一つのクラスと他の全てのクラスを識別する $K - 1$ 個の
2クラス識別関数 $f_j(\mathbf{x})(j = 1, \dots, K - 1)$ を用意し

$$\text{識別クラス} = \begin{cases} C_j & (\text{ある } j \text{ について } f_j(\mathbf{x}) > 0 \text{ の場合}) \\ C_K & (\text{すべての } j \neq K \text{ について} \\ & f_j(\mathbf{x}) < 0 \text{ の場合}) \end{cases}$$

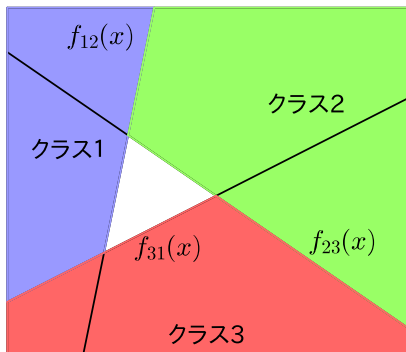


- 複数の識別関数が > 0 となる空白領域のクラスが決定できない
- 正のクラスの学習データ数が負のクラスの学習データ数に比べて極端に少なくなる

一対一 (one-versus-one)

クラス i と j を識別する $K(K-1)/2$ 個のクラス識別関数 $f_{ij}(\mathbf{x}) (1 \leq i < j \leq K)$ を用意し、 $K(K-1)/2$ 個の識別関数の多数決で識別クラスを決める。

$$\begin{cases} f_{ij}(\mathbf{x}) > 0 \Rightarrow C_i \text{ に一票} \\ f_{ij}(\mathbf{x}) < 0 \Rightarrow C_j \text{ に一票} \end{cases}$$



- 一対他と同じく空白領域が存在
- 手書き数字識別
10クラス → 識別関数 45 個
関係のある識別関数 9 個
過半数を取れない可能性あり

最大識別関数法 (1/2)

K 個の線形識別関数を用意し

$$\text{識別クラス} = \arg \max_j f_j(\mathbf{x}) = \arg \max_j (\mathbf{w}_j^T \mathbf{x} + w_{j0})$$

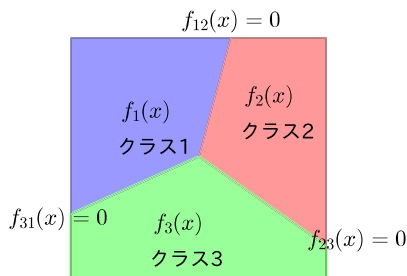
クラス i と j の識別境界は $f_i(\mathbf{x}) = f_j(\mathbf{x})$ となるので

$$f_{ij}(\mathbf{x}) = (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$$

この式を満たす $K - 1$ 個の識別境界ができる。

この識別境界は 2 クラスの場合の識別境界と同じ。

最大識別関数法 (2/2)



- 識別境界上を除けばどの点でも必ずどれかの識別関数が最大となる
 - クラスが決定できない領域が存在しない
- 各クラスの占める領域は単連結で凸となる

単連結

穴の開いていない領域の性質

凸

領域内の任意の2点を結ぶ直線上のすべての点が、その領域に含まれている時の性質

最小2乗誤差基準によるパラメータの推定(1/4)

係数ベクトル、 i 番目の学習用入力ベクトルはバイアスを含めて

$$\mathbf{w} = (w_0, w_1, \dots, w_d)^T$$

$$\mathbf{x}_i = (x_{i0} = 1, x_{i1}, \dots, x_{id})^T$$

線形識別関数

$$f(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_dx_d = \mathbf{w}^T \mathbf{x}$$

入力ベクトル \mathbf{x}_i が所属するクラスは、教師入力 t_i により

$$t_i = \begin{cases} +1 & (\mathbf{x}_i \in C_1) \\ -1 & (\mathbf{x}_i \in C_2) \end{cases}$$

のように与えられるものとする。

最小2乗誤差基準によるパラメータの推定 (2/4)

学習数を N とし、以下を定義.

$$\begin{aligned} \text{データ行列:} \quad \mathbf{X} &= (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \\ \text{教師ベクトル:} \quad \mathbf{t} &= (t_1, \dots, t_N)^T \end{aligned}$$

識別関数の出力値と教師入力の差を 2 乗誤差で評価

評価関数 $E(\mathbf{w})$

$$E(\mathbf{w}) = \sum_{i=1}^N (t_i - f(\mathbf{x}_i))^2 = (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})$$

$E(\mathbf{w})$ を最小にするパラメータ \mathbf{w} は

$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} &= -2\mathbf{X}^T (\mathbf{t} - \mathbf{X}\mathbf{w}) = 0 \\ \Rightarrow \hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \end{aligned}$$

最小2乗誤差基準によるパラメータの推定 (3/4)

正規方程式

$$\hat{\boldsymbol{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{t}$$

学習データに対する予測値

$$\hat{\boldsymbol{t}} = \mathbf{X} \hat{\boldsymbol{w}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{t}$$

ハット行列

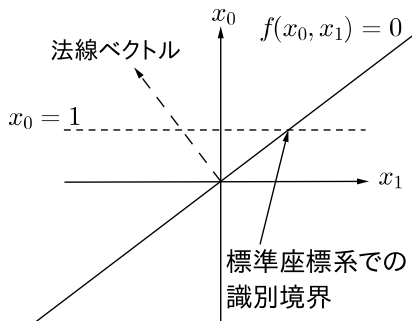
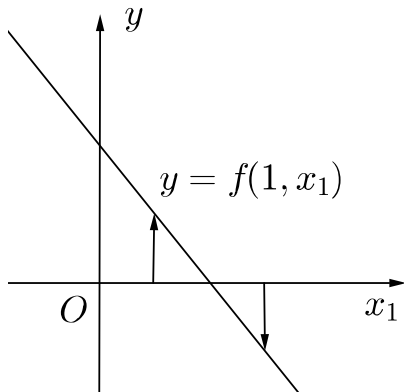
$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

教師データ \boldsymbol{t} を予測値 $\hat{\boldsymbol{t}}$ に変換する行列

最小2乗誤差基準によるパラメータの推定 (4/4)

標準座標系

x_0 を1に固定し、 $(x_1, f(1, x_1))$ 平面による表現



同次座標系

x_0 を固定せず独立した座標として扱うことで、2次元空間 (x_0, x_1) 内の識別関数として表現する

多クラス問題への拡張 (1/2)

最大識別関数法

識別関数

$$f_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} \quad (k = 1, \dots, K)$$

N 個の学習データ $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$

教師データの行列 $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_N)^T$, $\mathbf{t}_i = \underbrace{(0, \dots, 1, \dots, 0)^T}_{K \text{ 個の要素}}$

2 乗誤差を最小にするパラメータは

$$\widehat{\mathbf{W}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}$$

多クラス問題への拡張 (2/2)

識別関数

$$G(\mathbf{x}) = \widehat{\mathbf{W}}^T \mathbf{x} = (\mathbf{w}_1, \dots, \mathbf{w}_K)^T \mathbf{x} = (g_1(\mathbf{x}), \dots, g_K(\mathbf{x}))^T$$

識別規則

$$\text{識別クラス} = \arg \max_j g_j(\mathbf{x})$$

複数 (> 2) のクラスが一直線上に並んでいるような分布をしている場合、最大識別関数法ではうまく識別できない。

線形判別分析

- 線形識別関数: d 次元ベクトル \boldsymbol{x} を、ベクトル \boldsymbol{w} 上のスカラー関数 $f(\boldsymbol{x})$ に写像している
- 最小2乗誤差法: 教師データにできるだけ忠実になるよう線形識別関数を求めた

線形判別分析

1次元に写像されたとき、クラス間の分布ができるだけ重ならないような写像方向を見つける。

フィッシャーの線形判別関数 (1/2)

2クラス (C_1, C_2) 問題について

各クラスの学習データ数

$$N_1, N_2$$

線形識別関数

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

平均ベクトル

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i \in C_k} \mathbf{x}_i$$

写像

$$\boldsymbol{\mu}_k \mapsto m_k = \mathbf{w}^T \boldsymbol{\mu}_k$$

平均の差

$$m_1 - m_2 = \mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

平均の差が大きいほど、クラス分離が良くなる。

フィッシャーの線形判別関数 (2/2)

クラス間変動 平均の差の 2 乗 $(m_1 - m_2)$

クラス内変動 1次元に写像後の分散 $S_k^2 = \sum_{i \in C_k} (y_i - m_k)^2$

全クラス内変動 $S_1^2 + S_2^2$

フィッシャーの基準

クラス間変動とクラス内変動の比

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{S_1^2 + S_2^2}$$

を最大にする \mathbf{w} を見つけること.

多クラス問題への拡張

他クラス $K (> 2)$ の場合, $d (> K)$ 次元のデータをただか $K - 1$ 次元の特徴空間に写像する線形変換行列を見つける問題になる.

⇒ 識別境界は計算できない.

ロジスティック回帰

線形識別関数 $y = \mathbf{w}^T \mathbf{x}$ は、識別境界から離れるに従って、関数値の大きさが線形に上昇し続ける。

ロジスティック回帰

関数値を区間 $(0, 1)$ に制限して確率的な解釈を可能にする。

2クラス問題において、クラス C_1 の事後確率 $P(C_i|\mathbf{x})$ は

$$P(C_1|\mathbf{x}) = \frac{p(\mathbf{x}|C_1)P(C_1)}{p(\mathbf{x}|C_1)P(C_1) + p(\mathbf{x}|C_2)P(C_2)}$$

$$a = \ln \frac{p(\mathbf{x}|C_1)P(C_1)}{p(\mathbf{x}|C_2)P(C_2)}$$

とおけば

$$P(C_1|\mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

ロジスティック回帰

ロジスティック関数

$\sigma(a)$ をロジスティック関数(シグモイド関数)と呼び, $y = \sigma(x)$ は無限区間 $(-\infty, \infty)$ を区間 $(0, 1)$ に写像する圧縮関数.
 $\sigma(-a) = 1 - \sigma(a)$ のような対称性を示す.

ロジスティック関数の逆関数(ロジット関数)

$$a = \ln \left(\frac{\sigma(a)}{1 - \sigma(a)} \right) = \ln \frac{P(C_1|\mathbf{x})}{P(C_2|\mathbf{x})}$$

オッズ (事後確率の比) $\frac{P(C_1|\mathbf{x})}{P(C_2|\mathbf{x})}$

ログオッズ (オッズの対数) $\ln \frac{P(C_1|\mathbf{x})}{P(C_2|\mathbf{x})}$

ロジスティック回帰モデル

N 人の喫煙量 $\hat{\mathbf{x}} = (x_1, \dots, x_N)^T$ を観測したとき、肺がんになる確率を

$$p(1|x_1, \dots, x_N) = f(\mathbf{x}) = \frac{1}{1 + \exp(-(w_0 + w_1x_1 + \dots + w_Nx_N))}$$

で表す。($\hat{\mathbf{x}}$ にはバイアス項に対応する 1 を追加)

$$\begin{aligned} a &= \mathbf{w}^T \mathbf{x} \\ \Rightarrow f(\mathbf{x}) = \sigma(a) &= \frac{1}{1 + \exp(-a)} = \frac{\exp a}{1 + \exp a} \end{aligned}$$

$f(\mathbf{x})$ は非線形 \Rightarrow 学習データの線形関数である a をロジスティック関数で非線形変換したモデルで事象を表現していることになる。

ロジスティック回帰モデル

ロジスティック関数の逆関数であるロジット関数は,

$$a = \ln \frac{p(1|\mathbf{x})}{1 - p(1|\mathbf{x})} = \mathbf{w}^T \mathbf{x}$$

で, オッズは

$$\frac{p(1|\mathbf{x})}{1 - p(1|\mathbf{x})} = \frac{p(1|\mathbf{x})}{p(0|\mathbf{x})} = \exp(\mathbf{w}^T \mathbf{x})$$

で表される.

ロジスティック回帰モデル

x 中の x_1 が1 増えた状態 $\hat{x} = (1, (x_1 + 1), x_2, \dots, x_N)$ を考えると, x とのオッズ比は

$$\frac{\frac{p(1|\hat{x})}{1-p(1|\hat{x})}}{\frac{p(1|x)}{1-p(1|x)}} = \frac{\exp(w_0 + w_1(x_1 + 1) + w_2x_2 + \dots + w_Nx_N)}{\exp(w_0 + w_1x_1 + w_2x_2 + \dots + w_Nx_N)} = \exp w_1$$

x_1 が 1 単位分増えることによるオッズの増加分が $\exp w_1$

パラメータの最尤推定

2クラスロジスティック回帰モデルのパラメータの最尤推定を考える。

確率変数	t	:	モデルの出力
$t = 1$ となる確率	$P(t = 1)$	=	α
$t = 0$ となる確率	$P(t = 0)$	=	$1 - \alpha$

確率変数 t はパラメータ α をもつベルヌーイ試行

$$f(t|\alpha) = \alpha^t(1 - \alpha)^{1-t} \quad (t = 0 \text{ または } 1)$$

に従うので、 N 回の試行に基づく尤度関数は

$$L(\alpha_1, \dots, \alpha_N) = \prod_{i=1}^N f(t_i|\alpha_i) = \prod_{i=1}^N \alpha_i^{t_i} (1 - \alpha_i)^{(1-t_i)}$$

パラメータの最尤推定

負の対数尤度関数

$$\zeta(\alpha_1, \dots, \alpha_N) = - \sum_{i=1}^N (t_i \ln \alpha_i + (1 - t_i) \ln(1 - \alpha_i))$$

この評価関数は交差エントロピー型誤差関数とよばれる。
最尤推定とは、この交差エントロピー型誤差関数を最小にするパラメータ w を得ることである。

- 解析的に求められない
- 最急降下法やニュートン-ラフソン法などで数値的に解を求めることになる

多クラス問題への拡張

各クラスごとに線形変換

$$a_k = \mathbf{w}_k^T \mathbf{x} \quad (k = 1, \dots, K)$$

を求め、事後確率を

$$P(C_k | \mathbf{x}) = \pi_k(\mathbf{x}) = \frac{\exp a_k}{\sum_{j=1}^K \exp a_j}$$

で計算して、最大事後確率を与えるクラスに分類する。
この関数をソフトマップ関数という。

多クラス問題への拡張

負の対数尤度関数

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(T | \mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{i=1}^K \sum_{k=1}^K t_{ik} \pi_{ik}$$

となる.

各 \mathbf{w}_j の最尤推定は評価関数を \mathbf{w}_j で微分して 0 とおけば求められる. \Rightarrow 2クラス問題同様, 解析的には解けない.