

## 0 このゼミの概要

春季休業を利用して，学習理論の入門的な講義を 4 回行います．ベイズ推定の話を中心に以下のようなトピックスを扱います：

- 第 1 回 (2/23) 多項式による曲線近似
- 第 2 回 (3/02) 線形回帰
- 第 3 回 (3/09) 分類問題
- 第 4 回 (3/16) サポートベクタマシン入門

このゼミの教科書および参考書として以下のものを挙げておきます：

- 教科書: 平岡裕章, 応用数理 II 講義ノート (コピーを初回に配布予定)
- 参考書: Bishop, C. *Pattern Recognition and Machine Learning*, Springer, (2006)

また，簡単なレジюмеのようなものを作成して配布する予定です．

## 1 多項式による曲線近似

今回および次回は  $N$  個の既知データ  $(x_1, t_1), \dots, (x_N, t_N)$  があるときに，未知の  $\hat{x}$  に対する出力  $\hat{t}$  を予測する問題を考察する．

今回は，既知データ  $(x_i, t_i)$  を近似する多項式  $y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M$  の係数を推定して，未知の  $\hat{x}$  に対する出力  $\hat{t}$  を  $\hat{t} = y(\hat{x}, \mathbf{w})$  と予測することを考えよう．

### 1.1 最小二乗法

方針：エラー関数

$$\text{Err}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2$$

を最小にするような  $\mathbf{w}$  を求める．

解法：Err は  $\mathbf{w}$  に関して 2 次なので

$$\frac{\partial}{\partial \mathbf{w}} \text{Err}(\mathbf{w}) = \begin{pmatrix} \frac{\partial}{\partial w_0} \text{Err}(\mathbf{w}) \\ \frac{\partial}{\partial w_1} \text{Err}(\mathbf{w}) \\ \vdots \\ \frac{\partial}{\partial w_N} \text{Err}(\mathbf{w}) \end{pmatrix} = \mathbf{0}$$

を解くことによって  $\mathbf{w} = (w_0, w_1, \dots, w_N)'$  が推定できる．

問題点：過剰適合の問題がある (資料 1 参照)．

## 1.2 最尤推定法

ここではベイズ推定への準備としてデータが

$$t_i = y(x_i, \mathbf{w}) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \beta^{-1})$$

によって生成されていると仮定し, 最尤法によって $\mathbf{w}$ を推定することを考える.

注意: 出力 $t$ の確率密度関数は

$$p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x | \mathbf{w}), \beta^{-1}), \quad \text{where} \quad \mathcal{N}(t | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

である.

方針: 尤度関数

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t_1 | y(x_1, \mathbf{w}), \beta^{-1}) \cdots \mathcal{N}(t_N | y(x_N, \mathbf{w}), \beta^{-1})$$

を最大にする $\mathbf{w} = \mathbf{w}_{\text{ML}}$ を求める.

解法: 対数尤度関数が

$$\log p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi)$$

と変形できるので,  $\mathbf{w}_{\text{ML}}$ は $\text{Err}(\mathbf{w})$ を最小化にすることによって得られる.

## 1.3 ベイズ推定法

資料1によると, 過剰適合している $\mathbf{w}$ がとても大きいことが分かる. そこで事前に $\mathbf{w}$ はあまり大きな値はとらないという情報を与えて $\mathbf{w}$ を推定しよう.

方針: 事前確率を例えば $p(\mathbf{w}) = \mathcal{N}(w_0 | 0, \alpha^{-1}) \cdots \mathcal{N}(w_M | 0, \alpha^{-1})$ で導入し, 事後確率

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \beta) = \frac{p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w})}{p(\mathbf{t} | \mathbf{x}, \beta)}$$

を最大にする $\mathbf{w} = \mathbf{w}_{\text{MAP}}$ を求める.

解法: 事後確率は

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \beta) = C \exp \left\{ - \left( \frac{\beta}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\alpha}{2} \sum_{i=0}^M w_i^2 \right) \right\}$$

と変形できるので,

$$\widetilde{\text{Err}}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\alpha}{2\beta} \sum_{i=0}^M w_i^2$$

を最小にするような $\mathbf{w}$ を求めればよい.

## 2 線形回帰

前回の多項式による曲線近似を少し一般化して,

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}' \boldsymbol{\phi}(\mathbf{x})$$

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x})), \quad \phi_0(\mathbf{x}) = 1, \quad \mathbf{x} \in \mathbb{R}^D, \quad \mathbf{w} \in \mathbb{R}^M$$

とおいて, 入力 $\mathbf{x}$ に対し出力 $t$ が

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

で与えられているときの係数 $\mathbf{w}$ の推定を考える.

### 2.1 最尤推定法

方針: 対数尤度関数

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}' \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}), \quad \mathbf{t} = (t_1, \dots, t_N)', \quad \mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$$

の最大点 $\mathbf{w}_{\text{ML}}, \beta_{\text{ML}}$ を求める (以下では,  $\mathbf{X}$ を省略する.)

解法: 対数尤度関数は

$$\log p(\mathbf{t} | \mathbf{w}, \beta) = \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi) - \beta \text{Err}(\mathbf{w}), \quad \text{Err}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}' \boldsymbol{\phi}(\mathbf{x}_n))^2$$

と変形できるので $\mathbf{w}_{\text{ML}}, \beta_{\text{ML}}$ を求めるには

$$\frac{\partial}{\partial \mathbf{w}} \log p(\mathbf{t} | \mathbf{w}, \beta) = \mathbf{0}, \quad \frac{\partial}{\partial \beta} \log p(\mathbf{t} | \mathbf{w}, \beta) = 0$$

を解けばよい. これらはそれぞれ

$$\left( \sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)' \right) \mathbf{w} = \sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) t_n, \quad \frac{N}{2\beta} - \text{Err}(\mathbf{w}) = 0 \quad (\spadesuit)$$

と変形できる. ここで

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

を用いると $(\spadesuit)$ の前者は

$$\Phi' \Phi \mathbf{w} = \Phi' \mathbf{t}$$

と変形できるので, 結局 $\mathbf{w}_{\text{ML}}, \beta_{\text{ML}}$ は

$$\mathbf{w}_{\text{ML}} = (\Phi' \Phi)^{-1} \Phi' \mathbf{t}, \quad \beta_{\text{ML}} = \frac{N}{\sum_{n=1}^N (t_n - \mathbf{w}'_{\text{ML}} \boldsymbol{\phi}(\mathbf{x}_n))^2}$$

によって得られる.

## 2.2 ベイズ推定法

前回見たように，最尤推定法では過剰適合の問題がある．そこで前回同様に $w$ の事前分布を与えることによって，係数を推定する．

方針：事前確率を

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1}I)$$

で導入し，事後確率

$$p(\mathbf{w} | \mathbf{t}, \alpha, \beta) = c p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha)$$

を最大化する．

解法：事後確率は

$$p(\mathbf{w} | \mathbf{t}, \alpha, \beta) = c' \exp \varphi(\mathbf{w}), \quad \varphi(\mathbf{w}) = -\frac{\beta}{2} \sum_{n=1}^N (\mathbf{w}' \phi(\mathbf{x}_n) - t_n)^2 - \frac{\alpha}{2} \mathbf{w}' \mathbf{w}$$

と変形できる．さらに， $\varphi(\mathbf{w})$  は

$$\varphi(\mathbf{w}) = -\frac{1}{2}(\mathbf{w} - \mathbf{a})' B^{-1}(\mathbf{w} - \mathbf{a}), \quad B^{-1} = \alpha I + \beta \Phi' \Phi, \quad \mathbf{a} = \beta B \Phi' \mathbf{t}$$

と変形できるのでベイズ推定により推定値は

$$\mathbf{w}_{\text{MAP}} = \mathbf{a} = \beta B \Phi' \mathbf{t}$$

で与えられる．

## 2.3 時系列ベイズ推定

前節は  $N$  個全てのデータを用いて一度に推定する方法であった．ここでは  $N$  個のデータ  $(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)$  が順次手に入るとき，その都度  $w$  を学習していく方法を考える．

アイデア：事後確率を次のステップの事前確率に使う．

尤度関数	事後確率		
	$p(\mathbf{w})$		
$p(t_1   \mathbf{w})$	$\rightarrow$	$p(\mathbf{w}   t_1) = c_1 p(t_1   \mathbf{w}) p(\mathbf{w})$	ステップ 1
$p(t_2   \mathbf{w})$	$\rightarrow$	$p(\mathbf{w}   t_1, t_2) = c_2 p(t_2   \mathbf{w}) p(\mathbf{w}   t_1)$	ステップ 2
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$p(t_k   \mathbf{w})$	$\rightarrow$	$p(\mathbf{w}   t_1, \dots, t_k) = c_k p(t_k   \mathbf{w}) p(\mathbf{w}   t_1, \dots, t_{k-1})$	ステップ $k$

具体例については資料 2 を参照のこと．

### 3 分類問題

このゼミの後半の2回は入力された  $\boldsymbol{x} \in \mathbb{R}^D$  に対して, それをどちらかのクラス  $C_1, C_2$  に割り当てる分類問題を考察する.

方針: 既知データを用いて  $\mathbb{R}^D$  を  $\mathbb{R}^D = C_1 \cup C_2$  と分割する.

#### 3.1 非確率的手法を用いた解法

クラス数が2の分類問題は超平面  $y(\boldsymbol{x}) = \boldsymbol{w}'\boldsymbol{x} + w_0 = 0$  を1つ定め,  $y(\boldsymbol{x}) \geq 0$  のとき  $\boldsymbol{x} \in C_1$ ,  $y(\boldsymbol{x}) < 0$  のとき  $\boldsymbol{x} \in C_2$  と分類すればよい. しかしクラス数が3以上のときに拡張できるように次の分類方法を考える.

分類方法: 2個の線形関数

$$y_1(\boldsymbol{x}) = \boldsymbol{w}'_1\boldsymbol{x} + w_{10}, \quad y_2(\boldsymbol{x}) = \boldsymbol{w}'_2\boldsymbol{x} + w_{20}$$

を導入し,

$$y_1(\boldsymbol{x}) \geq y_2(\boldsymbol{x}) \Leftrightarrow \boldsymbol{x} \in C_1, \quad y_1(\boldsymbol{x}) < y_2(\boldsymbol{x}) \Leftrightarrow \boldsymbol{x} \in C_2$$

で分類分けを考える.

記号: 簡略化のため  $\boldsymbol{t} = \boldsymbol{t}(\boldsymbol{x})$  を,  $\boldsymbol{t} = (1, 0)'$   $\Leftrightarrow \boldsymbol{x} \in C_1$ ,  $\boldsymbol{t} = (0, 1)'$   $\Leftrightarrow \boldsymbol{x} \in C_2$  で定める.

即ち, 考える問題は既知データ  $(\boldsymbol{x}_1, \boldsymbol{t}_1), \dots, (\boldsymbol{x}_N, \boldsymbol{t}_N)$  を用いて

$$\widetilde{W} = \begin{pmatrix} w_{10} & w_{20} \\ \boldsymbol{w}_1 & \boldsymbol{w}_2 \end{pmatrix} = (\tilde{\boldsymbol{w}}_1 \quad \tilde{\boldsymbol{w}}_2) \in M_{D+1,2}(\mathbb{R})$$

を学習する問題になる. ここで

$$y_k > y_j \Leftrightarrow (y_k - 1)^2 + y_j^2 < y_k^2 + (y_j - 1)^2$$

が成り立つことを使うと, 任意の  $i$  番目のデータに対して

$$\sum_{j=1}^2 (y_j(\boldsymbol{x}_i) - t_{ij})^2 = (y_1(\boldsymbol{x}_i) - t_{i1})^2 + (y_2(\boldsymbol{x}_i) - t_{i2})^2$$

が最小になるような  $\widetilde{W}$  は既知データの正しい分類を与えることがわかる.

方針: エラー関数

$$\text{Err}(\widetilde{W}) = \frac{1}{2} \text{tr}\{(\tilde{X}\widetilde{W} - T)(\tilde{X}\widetilde{W} - T)'\} \left( = \frac{1}{2} \sum_{n=1}^N \sum_{j=1}^2 (y_j(\boldsymbol{x}_i) - t_{ij})^2 \right)$$
$$\tilde{X} = \begin{pmatrix} 1 & \cdots & 1 \\ \boldsymbol{x}_1 & \cdots & \boldsymbol{x}_N \end{pmatrix}' = (\tilde{\boldsymbol{x}}_1 \quad \cdots \quad \tilde{\boldsymbol{x}}_N)', \quad T = (\boldsymbol{t}_1 \quad \cdots \quad \boldsymbol{t}_N)'$$

を最小にする  $\widetilde{W}_*$  を求める.

注意: 行列  $\tilde{X}\widetilde{W}$  の  $(i, j)$  成分は  $y_j(\boldsymbol{x}_i)$  である.

解法：次の性質

$$\frac{\partial}{\partial A} \text{tr}(A'BA) = (B + B')A, \quad \frac{\partial}{\partial A} \text{tr}(A'B) = B$$

を用いると,  $\widetilde{W}_*$  は具体的に

$$\widetilde{W}_* = (\widetilde{X}'\widetilde{X})^{-1}\widetilde{X}'T$$

と書ける.

### 3.2 確率的手法を用いた解法

ここでは, 確率的手法を用いた分類を考える. 要するに  $p(C_1|\mathbf{x})$  を求めて,  $p(C_1|\mathbf{x}) = 1/2$  となる超曲面を  $\mathbb{R}^D$  内で定め,  $\mathbf{x} \in C_1 \Leftrightarrow p(C_1|\mathbf{x}) \geq 1/2$  によって分類を与える.

方針: 以下のような流れで分類分けを定める境界を求める.

(1) データの発生について

$$p(\mathbf{x}|C_k) = \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)'\Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\} \quad (k = 1, 2),$$
$$p(C_1) = \pi, \quad p(C_2) = 1 - \pi, \quad (0 \leq \pi \leq 1)$$

を仮定し既知データ  $(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)$  を用いてパラメータ  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma, \pi$  を学習する. ここで  $t_i$  は,  $t_i = 1 \Leftrightarrow \mathbf{x}_i \in C_1, t_i = 0 \Leftrightarrow \mathbf{x}_i \in C_2$  であるとする.

(2) ベイズの定理を用いて

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

を求める (このとき,  $p(C_2|\mathbf{x}) = 1 - p(C_1|\mathbf{x})$  である.)

(3)  $p(C_1|\mathbf{x}) = 1/2$  を定める超平面が  $C_1$  と  $C_2$  の境界として定まる.

解法: まず次に注意する:

$$\mathbf{x}_n \in C_1 \Rightarrow p(\mathbf{x}_n, t_n = 1) = \pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma),$$
$$\mathbf{x}_n \in C_2 \Rightarrow p(\mathbf{x}_n, t_n = 0) = (1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)$$

これより  $\mathbf{t} = (t_1, \dots, t_N)$ ,  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  とおいたとき, 尤度関数が

$$p(\mathbf{t}, \mathbf{X} | \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma, \pi) = \prod_{n=1}^N \{\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \Sigma)\}^{t_n} \{(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \Sigma)\}^{1-t_n}$$

となる. 対数尤度関数を各パラメータで微分することにより, 最大点がそれぞれ

$$\pi^* = \frac{1}{N} \sum_{n=1}^N t_n, \quad \boldsymbol{\mu}_1^* = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \quad \boldsymbol{\mu}_2^* = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n, \quad \Sigma^* = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

であることが分かる．ここで記号の意味は以下の通り．

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1^*)(\mathbf{x}_n - \boldsymbol{\mu}_1^*)', \quad N_1 = \sum_{n=1}^N t_n,$$

$$S_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2^*)(\mathbf{x}_n - \boldsymbol{\mu}_2^*)', \quad N_2 = \sum_{n=1}^N (1 - t_n)$$

従ってパラメータが推定できた．次に (♡) は

$$p(C_1|\mathbf{x}) = \frac{1}{1 + \exp(-a)} =: \sigma(a), \quad a = \log \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)}$$

と変形できる．ここで  $\sigma(\cdot)$  は logistic sigmoid 関数と呼ばれる単調増加関数である．これに  $p(\mathbf{x}|C_k), p(C_k)$  を代入して整理すると

$$a = \mathbf{w}'\mathbf{x} + w_0$$

の形になる．ここで

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

である．結局 (♡) は  $p(C_1|\mathbf{x}) = \sigma(\mathbf{w}'\mathbf{x} + w_0)$  となる．ところで  $\sigma(a) = 1/2$  のとき  $a = 0$  であるから

$$p(C_1|\mathbf{x}) = \sigma(\mathbf{w}'\mathbf{x} + w_0) = \frac{1}{2}$$

は  $\mathbf{w}'\mathbf{x} + w_0 = 0$  となる．この超平面が分類分けを定める境界となる．

### 3.3 [補] 超平面の幾何的考察

$y(\mathbf{x}) = \mathbf{w}'\mathbf{x} + w_0 = 0$  は  $\mathbb{R}^D$  の超平面を定める．これを  $P$  とおくと以下の性質が成り立つ：

- $\mathbf{x}_A, \mathbf{x}_B \in P$  のとき  $y(\mathbf{x}_A) = y(\mathbf{x}_B)$  なので  $\mathbf{w}'(\mathbf{x}_A - \mathbf{x}_B) = 0$  となる．従って  $\mathbf{w}$  は  $P$  に直交するベクトルである．
- 原点  $O$  から超平面  $P$  までの最短距離  $r_n$  は， $\mathbf{x} \in P$  をとると次のように表される：

$$r_n = \frac{\mathbf{w}'\mathbf{x}}{\|\mathbf{w}\|} = \frac{-w_0}{\|\mathbf{w}\|}$$

- $\mathbf{x} \in \mathbb{R}^D$  に対して， $\mathbf{x}_\perp \in P$  を  $\mathbf{x}$  の  $P$  への射影とすると，

$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}, \quad r \in \mathbb{R}$$

と書ける．このとき  $r$  は次のように求められる．

$$y(\mathbf{x}) = \mathbf{w}'\mathbf{x}_\perp + r \frac{\mathbf{w}'\mathbf{w}}{\|\mathbf{w}\|} + w_0 = r\|\mathbf{w}\|, \quad r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

これらの性質はサポートベクタマシンのときに使用する．

## 4 サポートベクタマシン入門

今回は既知データから超平面

$$y(\mathbf{x}) = \mathbf{w}'\mathbf{x} + w_0 = 0$$

を学習させてデータの分類を行った。今回は少し一般化して

$$y(\mathbf{x}) = \mathbf{w}'\phi(\mathbf{x}) + w_0, \quad \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_D(\mathbf{x}))' \quad (\diamond)$$

の形で分類問題を考える。このとき制約条件つき極値問題に帰着させるので、ラグランジュ未定乗数法について触れておく。

### 4.1 [補] ラグランジュ未定乗数法

ここでは2種類の制約条件つき極値問題を扱う。

- (1) 制約条件:  $g(\mathbf{x}) = 0$ ,  $\mathbf{x} \in \mathbb{R}^D$  の下で関数  $f(\mathbf{x})$  を最大にする点  $\mathbf{x}^*$  を求める。

補題:  $g(\mathbf{x}) = 0$  が定める超曲面を  $S$  とする。このとき  $\mathbf{x} \in S$  に対し

$$\nabla g(\mathbf{x}) = \left( \frac{\partial g}{\partial x_1}, \dots, \frac{\partial g}{\partial x_D} \right)$$

は  $S$  に直交する。

補題:  $f(\mathbf{x})$  が  $\mathbf{x}^* \in S$  で極大になるなら,  $\nabla f(\mathbf{x}^*)$  は  $S$  に直交する。

上の2つの補題から,  $\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = \mathbf{0}$  となる  $0 \neq \lambda \in \mathbb{R}$  が存在する。

従って, ラグランジュ関数を

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

で定義すると, 問題 (1) が極大値を持つ必要条件是

$$\nabla_{\mathbf{x}} L = \mathbf{0}, \quad \frac{\partial}{\partial \lambda} L = 0$$

である。これを満たす  $(\mathbf{x}^*, \lambda^*)$  を見つける方法をラグランジュ未定乗数法という。

- (2) 制約条件:  $g(\mathbf{x}) \geq 0$ ,  $\mathbf{x} \in \mathbb{R}^D$  の下で関数  $f(\mathbf{x})$  を最大にする点  $\mathbf{x}^*$  を求める。

制約条件は  $g(\mathbf{x}^*) > 0$  か  $g(\mathbf{x}^*) = 0$  かで場合分けをする。前者は無条件の極値問題であり, 後者は (1) と同じ条件付き極値問題となる。従ってラグランジュ関数

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

(但し  $g(\mathbf{x}) \geq 0$ ,  $\lambda \geq 0$ ,  $\lambda g(\mathbf{x}) = 0$ ) の極値を与える  $(\mathbf{x}^*, \lambda^*)$  を探す方法が (2) に対するラグランジュ未定係数法である。

- (3) 制約条件が複数:  $g_j(\mathbf{x}) = 0$  ( $j = 1, \dots, J$ ),  $h_k(\mathbf{x}) \geq 0$  ( $k = 1, \dots, K$ ) のとき

$$L(\mathbf{x}, \{\lambda_j\}, \{\mu_k\}) = f(\mathbf{x}) + \sum_{j=1}^J \lambda_j g_j(\mathbf{x}) + \sum_{k=1}^K \mu_k h_k(\mathbf{x})$$

(但し  $\mu_k \geq 0$ ,  $\mu_k h_k(\mathbf{x}) = 0$ ) とすればよい。



## 4.2 サポートベクタマシン

既知データを  $(\mathbf{x}_n, t_n), n = 1, \dots, N$  とする．ここで

$$t_n = \begin{cases} 1 & y(\mathbf{x}_n) > 0 \Leftrightarrow \mathbf{x}_n \in C_1 \\ -1 & y(\mathbf{x}_n) < 0 \Leftrightarrow \mathbf{x}_n \in C_2 \end{cases}$$

である．サポートベクタマシンはデータを分類する曲面  $y(\mathbf{x}) = 0$  のうち，最も近い既知データとの距離が最大になる曲面 を求める方法である．これは §3.3 と  $|y(\mathbf{x}_n)| = t_n y(\mathbf{x}_n)$  より

$$\arg \max_{\mathbf{w}, w_0} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}' \phi(\mathbf{x}_n) + w_0)] \right\}$$

で与えられる．ここで  $\{\cdot\}$  は  $\mathbf{w} \rightarrow \kappa \mathbf{w}, w_0 \rightarrow \kappa w_0 (0 \neq \kappa \in \mathbb{R})$  で不変なので，曲面に最も近い点  $\mathbf{x}_n$  で

$$t_n (\mathbf{w}' \phi(\mathbf{x}_n) + w_0) = 1$$

と規格化すると，次の制約条件付き極値問題になる：

制約条件： $t_i (\mathbf{w}' \phi(\mathbf{x}_i) + w_0) \geq 1 (i = 1, \dots, N)$  のもとで

$$\arg \max_{\mathbf{w}, w_0} \frac{1}{\|\mathbf{w}\|} = \arg \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2$$

を求める．

ラグランジュ関数を次で定義する：

$$L(\mathbf{w}, w_0, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}' \phi(\mathbf{x}_n) + w_0) - 1\}$$

$a_n \geq 0, t_n y(\mathbf{x}_n) \geq 1, a_n (t_n y(\mathbf{x}_n) - 1) = 0$  の下で  $L(\mathbf{w}, w_0, \mathbf{a})$  の極値を探せばよい．ここで

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n) = \mathbf{0}, \quad \frac{\partial L}{\partial w_0} = - \sum_{n=1}^N a_n t_n = 0$$

これより  $\mathbf{w}^* = \sum_{n=1}^N a_n t_n \phi(\mathbf{x}_n)$  なので， $L(\mathbf{w}, w_0, \mathbf{a})$  に代入すると次のようになる：

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m), \quad k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})' \phi(\mathbf{x}')$$

$a_n \geq 0, \sum_{n=1}^N a_n t_n = 0$  の下で極値  $\tilde{L}(\mathbf{a}^*)$  が求まれば， $\mathbf{w}^* = \sum_{n=1}^N a_n^* t_n \phi(\mathbf{x}_n)$  となる．これは MATLAB などの数式ソフトウェアで数値解を求めることができる．

学習解  $\mathbf{w}^*$  の性質： $(\diamond)$  に  $\mathbf{w}^*$  を代入し，さらに  $a_n (t_n y(\mathbf{x}_n) - 1) = 0$  に注意すると

$$\begin{aligned} y(\mathbf{x}) &= \mathbf{w}^{*'} \phi(\mathbf{x}) + w_0 = \sum_{n=1}^N a_n t_n k(\mathbf{x}_n, \mathbf{x}) + w_0 \\ &= \sum_{n \in S} a_n t_n k(\mathbf{x}_n, \mathbf{x}) + w_0, \quad \text{where } S = \{n \mid t_n y(\mathbf{x}_n) = 1\} \end{aligned}$$

となる． $n \in S$  となる  $\mathbf{x}_n$  は分類曲面  $y(\mathbf{x}) = 0$  の最も近くにある既知データなので未知データの分類の計算量が大幅に減少する．この最も近い点のことをサポートベクトルという．サポートベクトルのみを使ってデータの分類を行うのがサポートベクタマシンの特徴である．