

応用数理 II でやりたいこと

・学習理論の紹介

(線形回帰, 分類問題  
サポートベクターマシン (SVM))

・バイオインフォマクスへの学習理論の応用

参考書

・ Bishop, Pattern Recognition and Machine Learning, Springer

・ 阿久津 達也, バイオインフォマクスの数理とアルゴリズム, 共立出版

総科 C819, 6482

例題: 多項式による曲線近似

$N$  コの既知のデータ  $(x_1, t_1), \dots, (x_N, t_N)$

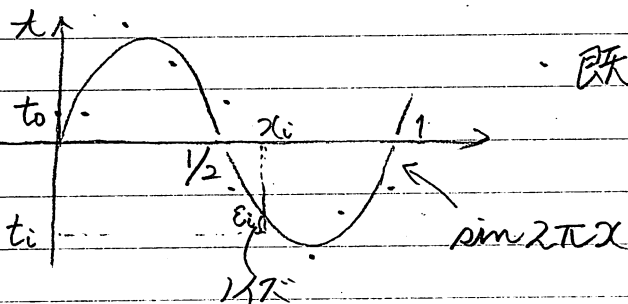
$x_i \in [0, 1], t_i \in \mathbb{R}$

の特徴を学習して未知のパラメータ  $\alpha$  に対する出力  $\hat{t}$  を予想せよ.

( $t_i$  は  $\sin 2\pi x_i + \underbrace{\varepsilon_i}_{\text{ノイズ}}$  によって生成) ... ☆

注意) 現実の問題では関係 ☆ は不明

(理解を助ける為に明示している)



多項式:  $y(x, w) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M$   
 で  $(x_i, t_i) \quad i=1, \dots, N$  を近似してあげよう。つまり  
 $w_0, w_1, \dots, w_M$  を学習(推定)してあげよう。

解法1  $y(x_i, w) = w_0 + w_1 x_i + \dots + w_M x_i^M = t_i \quad i=1, \dots, N$

$$\begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^M \\ 1 & x_2 & x_2^2 & \dots & x_2^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \dots & x_N^M \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_M \end{pmatrix} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix}$$

既知

未知

既知

$N=M$  の時以外は この方法は使えない  
 実用的ではない!

解法2. エラー関数

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

測定値と多項式のズレ  $E(w)$

を最小化する。

- $E(w)$  は  $w$  に関して 2次関数 (下に凸)
- $\frac{\partial E}{\partial w_i}$  は  $w$  に関して 1次関数

$$\Rightarrow \begin{pmatrix} \frac{\partial E}{\partial w_0}(w) \\ \vdots \\ \frac{\partial E}{\partial w_M}(w) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \quad M+1 \text{コ} \quad \dots (\$)$$

が定める  $M+1$  コの連立1次方程式をとくことで  
 $w = (w_0, \dots, w_M)$  が決定できる。

レポート1. (\$) を具体的に表しましょう ( $Aw = b$ )

問題点, (解法2の): OVER FITTING

学習理論ではどのようにして Over fitting を回避しているか？

⇒ 確率の導入 (ベイズ推定)

### 確率の復習

・ 離散確率変数  $x$

$x = \{x_1, \dots, x_n\}$   $x$  がとりうる事象

$P(x=x_i)$ :  $x=x_i$  とはる確率

$$\left(\sum_{i=1}^n P(x_i) = 1\right)$$

例. コイン投げ

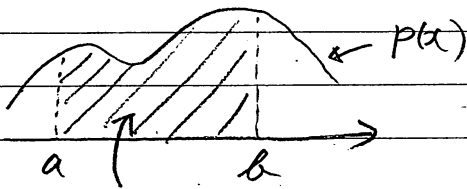
$x = \{\text{表}, \text{裏}\}$

$$P(\text{表}) = P(\text{裏}) = \frac{1}{2}$$

・ 連続確率変数  $x$

確率密度関数  $0 \leq p(x)$  s.t.  $\int_{-\infty}^{\infty} p(x) dx$

$$P(a \leq x \leq b) = \int_a^b f(x) dx$$



面積が確率に対応

例. 正規分布 (ガウス分布)

$$N(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

パラメータ

・  $x, y$  を確率変数としたとき

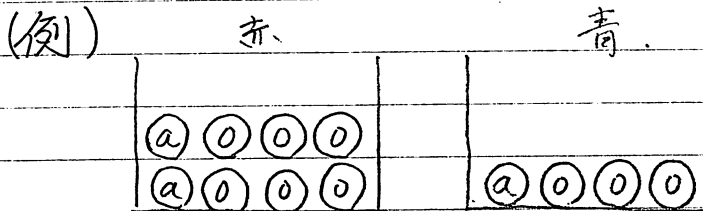
$p(x, y)$ : 結合確率

↑  $x$  かつ  $y$  をとる確率

$p(x|y)$ : 条件付き確率

( $y$  が既知の下での  $x$  の確率)

重要な性質  $P(x) = \sum_y P(x, y)$  ① 周辺化  
 $P(x, y) = P(x|y)P(y)$  ②



①: リンゴ, ②: オレンジ

ルール

赤い箱は 40% } の確率で選ばれる

青い箱は 60% }

各箱のフルーツは等確率で選ばれる

|   | ①  | ②  |
|---|--|--|
| 赤 | $\frac{4}{10} \times \frac{2}{4} = \frac{1}{5}$  | $\frac{4}{10} \times \frac{2}{4} = \frac{1}{5}$  |
| 青 | $\frac{6}{10} \times \frac{2}{4} = \frac{3}{10}$ | $\frac{6}{10} \times \frac{2}{4} = \frac{3}{10}$ |

$$\sum_{z \in Z} P(z) = 1$$

全対象

$$\Rightarrow \frac{1}{5} + \frac{1}{5} + \frac{3}{10} + \frac{3}{10} = 1$$

確率の表

$$P(\text{赤}) = \frac{4}{10}, P(\text{赤}, ①) + P(\text{赤}, ②) = \frac{1}{5} + \frac{1}{5} = \frac{4}{10}$$

$\Rightarrow$  ① が成り立っている

$$P(\text{赤}, ①) = \frac{1}{5} \quad P(①|\text{赤})P(\text{赤}) = \frac{2}{4} \times \frac{4}{10} = \frac{1}{5}$$

$\Rightarrow$  ② が成り立っている

ベイズの法則

$$P(x, y) = P(x|y)P(y) = P(y|x)P(x) \quad \text{よ}$$

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \quad \text{: ベイズの法則}$$

$x$ : 観測可能なデータ } とみなす

$y$ : 調べたい確率変数 }

$P(y)$ : 事前確率,  $P(y|x)$ : 事後確率

$P(x|y)$ : 尤度関数 ( $y$  についての関数)

$P(x)$ : 正規化因子  
(規格化)

$P(y|x)$  は  $y$  についての確率なので

$\sum_y P(y|x) = 1$  となる必要がある

$\sum_y c P(x|y) P(y) = 1 \Rightarrow c$  が一意に定まる  
 $\frac{1}{c} = P(x)$

(例の続き) ベイズの法則を体感しましょう

$P(\text{赤}) = \frac{6}{10}$ ,  $P(\text{青}) = \frac{4}{10}$  ← 事前確率

$P(a|\text{赤}) = \frac{1}{4}$   $P(o|\text{赤}) = \frac{3}{4}$

$P(o) = P(\text{赤}, o) + P(\text{青}, o) = \frac{3}{10} + \frac{3}{20} = \frac{9}{20}$

向 とり出されたフルーツがオレンジだとした

どちらの箱から選ばれたんだろう? → 観測データ

$$P(\text{赤}|o) = \frac{P(o|\text{赤})P(\text{赤})}{P(o)} = \frac{\frac{1}{4} \times \frac{6}{10} \times \frac{20}{9}}{\frac{9}{20}} = \frac{2}{3}$$

$$P(\text{青}|o) = 1 - P(\text{赤}|o) = \frac{1}{3}$$

## 期待値と分散

関数  $f(x)$  の期待値

$$E(f) = \begin{cases} \sum_x p(x) f(x) & \text{離散} \\ \int p(x) f(x) dx & \text{連続} \end{cases}$$

多変数の場合 ( $p(x, y)$ )

$$E_x(f) = \sum_x p(x) f(x, y) \quad x \text{ に関する期待値}$$

$$E_x(f|y) = \sum_x p(x|y) f(x, y) \quad x \text{ に関する}$$

条件付期待値

関数  $f(x)$  の分散

$$\text{var}(f) = E((f(x) - E(f))^2)$$

$$= E(f(x)^2) - 2E(f)f(x) + E(f)^2$$

Eの線形性

$$= E(f(x)^2) - 2E(f)E(f(x)) + E(f)^2$$

$$= E(f^2) - 2E(f)^2 + E(f)^2$$

$$= E(f^2) - E(f)^2$$

特 $\kappa$   $x$ の分散:  $\text{var}(x) = E((x - E(x))^2)$   
 $= E(x^2) - E(x)^2$

確率変数  $x, y$  の共分散

$$\text{var}(x, y) = E_{xy}((x - E_x(x))(y - E_y(y)))$$

$$= E_{xy}(xy) - E_x(x)E_y(y)$$

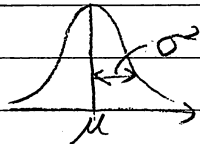
確率変数  $x, y$  の共分散

$$\text{var}(x, y) = E_{xy}((x - E_x(x))(y - E_y(y)))'$$

( $m, n$ ) 行列

がうす分布

$$N(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$



$$E(x) = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) dx = \mu$$

$$\text{var}(x) = \sigma^2$$

レポート  
(示せ)

この意味で  $\mu$  を平均,  $\sigma^2$  を分散と呼ぶ。

がうす分布に対する最大推定法

がうす分布:

$$N(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

$\mu, \sigma^2$  は未知である  $N(x|\mu, \sigma^2)$  から  $N$  個の独立なデータ

$$x = (x_1, \dots, x_N)$$

が得られたとする。これから  $\mu, \sigma^2$  を推定しよう!

$$P(x|\mu, \sigma^2) = N(x|\mu, \sigma^2) \cdots N(x_N|\mu, \sigma^2) \quad (\text{独立だから})$$

### 最大推定法の方針

: 尤度関数  $P(x|\mu, \sigma^2)$  の最大点,  $(\mu_{ML}, \sigma_{ML}^2)$  を探す  
Maximum Likelihood.

$$\begin{aligned} \ln P(x|\mu, \sigma^2) &= \sum_{n=1}^N \ln N(x_n|\mu, \sigma^2) \\ &= \sum_{n=1}^N \left[ -\ln(2\pi\sigma^2)^{1/2} - \frac{1}{2\sigma^2}(x_n - \mu)^2 \right] \\ &= -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \cdots \textcircled{1} \end{aligned}$$

①の最大点,  $\mu_{ML}, \sigma_{ML}^2$  を求めればよい。

$$\frac{d}{d\mu} \ln P(x|\mu, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^N (x_n - \mu) = 0$$

$$\Leftrightarrow \sum_{n=1}^N x_n - N\mu = 0$$

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \quad (\text{データの平均})$$

同様  $\sigma$

$$\frac{d}{d\sigma} \ln P(x|\mu, \sigma^2) = \frac{1}{2\sigma^4} \sum_{n=1}^N (x_n - \mu_{ML})^2 - \frac{N}{2\sigma^3} = 0$$

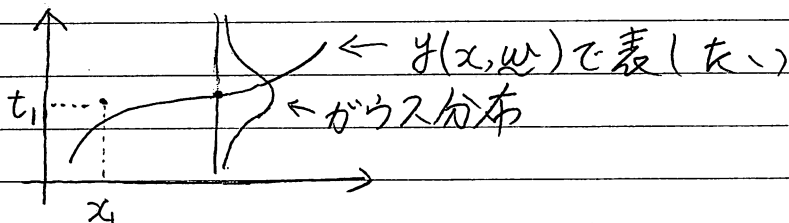
$$\Leftrightarrow \sum_{n=1}^N (x_n - \mu_{ML})^2 = N\sigma^2$$

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2 \quad (\text{データの分散})$$

これをもとに  $K$  多項式  $K$  による曲線近似 (オ1回) の例を再考

解法2の問題点,  $\Rightarrow$  over fitting ( $|w| \gg 1$ )

「解法2 = 最大推定法, over fitting  $\Leftarrow$  バイス推定」



$$t_i = y(x_i, w) + \epsilon_i \quad \epsilon_i \text{ をガウス分布でモデリングする}$$

$$P(t|x, w, \beta) = N(t|y(x|w), \beta^{-1})$$

$K$  上,  $T$  データが生成されているとする。

各データは独立なので

$$P(t|x, \omega, \beta) = N(t_1 | y(x_1, \omega), \beta_1^{-1}) \cdots N(t_N | y(x_N, \omega), \beta_N^{-1})$$

データに対する尤度関数

最大推定を用いて  $\omega$  を推定してみる。①と同様

$$\ln P(t|x, \omega, \beta) = \sum_{n=1}^N \ln N(t_n | y(x_n, \omega), \beta^{-1})$$

$$= -\frac{\beta^2}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

$$\omega_{ML} \text{ は } E(\omega) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2$$

を最小化することによって得られる。

ベイズ推定を導入

方針: 事後確率を最大化する。

事前確率  $P(\omega)$  を例えば次で導入する

$$P(\omega) = N(\omega_0 | 0, \alpha^{-1}) \cdots N(\omega_M | 0, \alpha^{-1})$$

(ベイズの法則より)

$$P(\omega|x, t, \beta) = \frac{P(t|x, \omega, \beta) P(\omega)}{P(t|x, \beta)} \leftarrow \text{規格化因子}$$

$$\propto P(t|x, \omega, \beta) P(\omega)$$

この事後確率を最大にする  $\omega$  を求めろ!

$$P(\omega|x, t, \beta) = C P(t|x, \omega, \beta) P(\omega)$$

$$= C \prod_{n=1}^N N(t_n | y(x_n, \omega), \beta^{-1}) \prod_{i=0}^M N(\omega_i | 0, \alpha^{-1})$$

$$= C \prod_{n=1}^N \frac{\beta^2}{(2\pi)^{1/2}} \exp\left\{-\frac{\beta^2}{2} (y(x_n, \omega) - t_n)^2\right\}$$

$$\times \prod_{i=0}^M \frac{\alpha^2}{(2\pi)^{1/2}} \exp\left(-\frac{\alpha}{2} \omega_i^2\right)$$

$$= C' C \exp\left\{-\left(\frac{\beta^2}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{\alpha}{2} \sum_{i=0}^M \omega_i^2\right)\right\}$$

( $C'$ :  $\omega$  に依存しない定数)

$\Rightarrow \frac{\beta^2}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{\alpha}{2} \sum_{i=0}^M \omega_i^2$  を最小にする  $\omega$  を探せばよい。

$\Leftrightarrow \tilde{E}(\omega) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \omega) - t_n)^2 + \frac{\alpha}{2\beta^2} \sum_{i=0}^M \omega_i^2$  を "  
 $E(\omega)$  (  $\omega$  が大きくなるのを防ぐ項  
(regularization 項)

$\Rightarrow$  over fitting を防げる。



まとめると

最大推定

over fitting が生じる

$E(w)$  の最小点

ベイズ推定

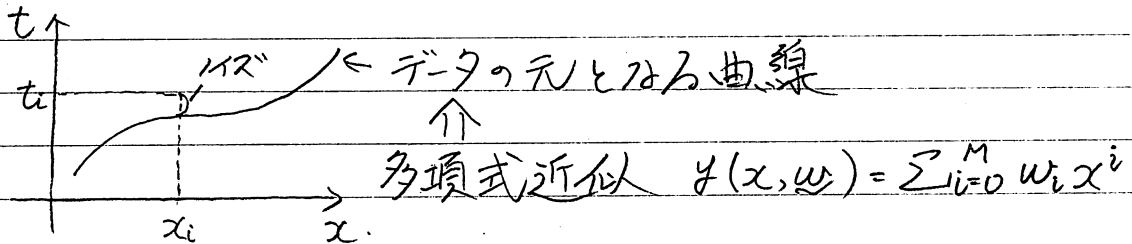
over fitting は抑制できる

$\tilde{E}(w)$  の最小点

今週と来週の内容: 線形回帰  
(linear Regression)

多項式による曲線(ノイズ)近似の一般化

前回まで



key は  $w$  に関して線形計算のみ

$$y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x) \in \mathbb{R} \quad \dots \textcircled{1}$$

$$= w' \phi(x)$$

ここで  $\phi(x) = (\phi_0(x), \dots, \phi_{M-1}(x))$

$\phi_0(x) = 1, x \in \mathbb{R}^D, w \in \mathbb{R}^M$

と一般化して考える。

つまり近似を行う基底を

$\{x^i\} \rightarrow \{\phi_j(x)\}$

と一般化

線形回帰問題: 入力  $x$  に対して出力  $t$  がノイズを  
供して観測されるとき

$$t = y(x, w) + \varepsilon$$

の形で対応関係をモデル化すること。

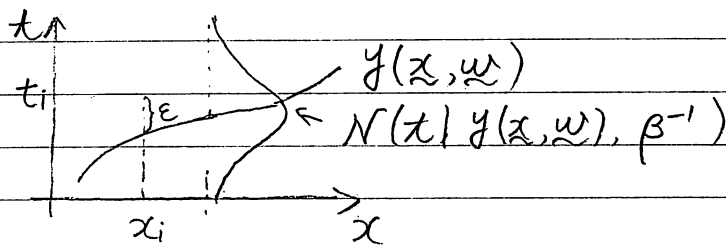
$\phi_0(x), \dots, \phi_{M-1}(x)$ : 基底関数

注意① 多項式による曲線近似は

$D=1$ ,  $\phi_j(x) = x^j$  に対応

② 基底関数として多項式以外では  
 ガウス関数, 三角関数 等々

前回同様,  $\varepsilon$  をガウス分布でモデル化する



$$p(t | x, w, \beta) = N(t | y(x, w), \beta^{-1})$$

既知データ:  $(x_1, t_1), \dots, (x_N, t_N)$

を用いて,  $w, \beta$  を学習したい.

各データは独立  $\Rightarrow$  データに対する尤度関数は

$$p(t | X, w, \beta) = N(t_1 | y(x_1, w), \beta^{-1}) \times \dots \times N(t_N | y(x_N, w), \beta^{-1}) \\ = N(t_1 | w' \phi(x_1), \beta^{-1}) \times \dots \times N(t_N | w' \phi(x_N), \beta^{-1})$$

ここで  $t = (t_1, \dots, t_N)$ ,  $X = (x_1, \dots, x_N)$

解法1: 最大推定法

$\Leftrightarrow p(t | X, w, \beta)$  の最大点,  $w_{ML}, \beta_{ML}$  を求める

$\Leftrightarrow \log p(t | X, w, \beta)$  の " "

以後  $X$  は省略

$$\log p(t | w, \beta) = \sum_{n=1}^N \log N(t_n | w' \phi(x_n), \beta^{-1})$$

$$= \sum_{n=1}^N \log \frac{\beta^{1/2}}{(2\pi)^{1/2}} \exp \left\{ -\frac{\beta}{2} (t_n - w' \phi(x_n))^2 \right\}$$

$$= \frac{N}{2} \log \beta - \frac{N}{2} \log(2\pi) - \beta E(w)$$

$$\text{ここで } E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w' \phi(x_n))^2$$

$w_{ML}$  を求める.

$$\frac{\partial}{\partial w_i} \log p(t | w, \beta) = 0 \Leftrightarrow \frac{\partial}{\partial w_i} E(w) = 0, \quad i=0, \dots, M-1$$

$$f') \sum_{n=1}^N (t_n - \underline{w}' \phi(x_n)) \phi_i(x_n) = 0 \quad i=0, \dots, M-1$$

$$\Rightarrow \sum_{n=1}^N t_n \phi_i(x_n) - \underline{w}' \sum_{n=1}^N \phi(x_n) \phi_i(x_n) = 0$$

横ベクトル表示

$$i=0, \dots, M-1$$

$$\Rightarrow \sum_{n=1}^N t_n \phi(x_n)' - \underline{w}' \sum_{n=1}^N \phi(x_n) \phi(x_n)' = 0$$

$$\Rightarrow \left( \sum_{n=1}^N \phi(x_n) \phi(x_n)' \right) \underline{w} = \sum_{n=1}^N \phi(x_n) t_n \quad (\$)$$

$M \times M$  行列

よって

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \dots & \phi_{M-1}(x_1) \\ \vdots & \ddots & \vdots \\ \phi_0(x_N) & \dots & \phi_{M-1}(x_N) \end{pmatrix} \quad \text{としとき}$$

$$\left. \begin{aligned} \cdot \Phi' \Phi &= \sum_{n=1}^N \phi(x_n) \phi(x_n)' \\ \cdot \Phi' \underline{t} &= \sum_{n=1}^N \phi(x_n) t_n \end{aligned} \right\} (\star)$$

を示しましょう

(\\$) と (\star) より

$$\Phi' \Phi \underline{w}_{ML} = \Phi' \underline{t}$$

$$\Rightarrow \underline{w}_{ML} = (\Phi' \Phi)^{-1} \Phi' \underline{t}$$

◦  $\beta_{ML}$  を求める.

$$\frac{\partial}{\partial \beta} \log P(\underline{t} | \underline{w}, \beta) = 0$$

$$\Rightarrow \sum_{n=1}^N \frac{\partial}{\partial \beta} \log p(t_n | \underline{w}, \beta) = 0$$

$$\frac{1}{\beta_{ML}} = \frac{1}{N} E(\underline{w}_{ML}) = \frac{1}{N} \sum_{n=1}^N (t_n - \underline{w}' \phi(x_n))^2 //$$

$x \in \mathbb{R}^D$  に対するガウス分布

$$N(x | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)' \Sigma^{-1} (x - \mu) \right\}$$

$\mu \in \mathbb{R}^D$ ,  $\Sigma: D \times D$  行列

平均 共分散行列

( $D=1$  での通常のガウス分布)

前回同様  $X$  は省略

## ベイズ推定法

事前確率を導入:

単位行列

★  $p(\underline{w}|\alpha) = N(\underline{w}|\underline{0}, \alpha^{-1}I)$  ... 多次元ガウス分布  
( $\underline{w}$ がベクトルだから)

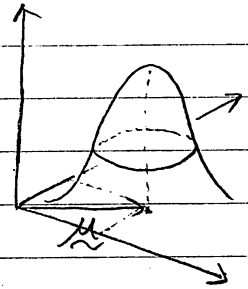
一般に  $\underline{x} \in \mathbb{R}^D$  に対する多次元ガウス分布は

$$N(\underline{x}|\underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\underline{x}-\underline{\mu})'\Sigma^{-1}(\underline{x}-\underline{\mu})\right\}$$

$\underline{\mu} \in \mathbb{R}^D$ .  $\Sigma: D \times D$  行列

平均 共分散行列

で与えられる. 当然,  $D=1$  では通常  
のガウス分布



実は  $\underline{\mu} = \underline{0}$ ,  $\Sigma = \alpha^{-1}I$  より

$$p(\underline{w}|\alpha) = \prod_{i=0}^{M-1} \frac{\alpha^{1/2}}{(2\pi)^{1/2}} \exp\left\{-\frac{\alpha}{2} w_i^2\right\}$$

これを用いてベイズの法則から事後確率を求めら:

$$p(\underline{w}|\underline{t}, \alpha, \beta) = c \underbrace{p(\underline{t}|\underline{w}, \beta)}_{\text{規格化定数}} p(\underline{w}|\alpha)$$

$$= c \prod_{n=1}^N \frac{\beta^{1/2}}{(2\pi)^{1/2}} \exp\left\{-\frac{\beta}{2} (\underline{w}'\phi(\underline{x}_n) - t_n)^2\right\}$$

$$\times \prod_{i=0}^{M-1} \frac{\alpha^{1/2}}{(2\pi)^{1/2}} \exp\left\{-\frac{\alpha}{2} w_i^2\right\}$$

$$= c' \exp\left\{-\frac{\beta}{2} \sum_{n=1}^N (\underline{w}'\phi(\underline{x}_n) - t_n)^2 - \frac{\alpha}{2} \underline{w}'\underline{w}\right\}$$

$$= c' \exp \varphi(\underline{w})$$

$$\varphi(\underline{w}) = -\frac{\beta}{2} \sum_{n=1}^N (\underline{w}'\phi(\underline{x}_n) - t_n)^2 - \frac{\alpha}{2} \underline{w}'\underline{w}$$

$$= -\frac{1}{2} (\underline{w} - \underline{a})' B^{-1} (\underline{w} - \underline{a})$$

と置いて  $\underline{a}, B$  を求めてみよう

(展開すると)

$$= -\frac{1}{2} (\underline{w}'B^{-1} - \underline{a}'B^{-1}) (\underline{w} - \underline{a})$$

$$= -\frac{1}{2} (\underline{w}'B^{-1}\underline{w} - \underline{a}'B^{-1}\underline{w} - \underline{w}'B^{-1}\underline{a} + \underline{a}'B^{-1}\underline{a})$$

⇒ ①  $\underline{w}$  の 2 次の係数から  $B^{-1}$  を求めら.

②  $\underline{w}$  の 1 次の係数と  $B^{-1}$  より  $\underline{a}$  を求めら

①  $\varphi(\underline{w})$  の  $\underline{w}$  に関して 2 次の項を比較すると

$$-\frac{1}{2} (\beta \sum_{n=1}^N (\underline{w}'\phi(\underline{x}_n))^2 + \alpha \underline{w}'\underline{w}) = -\frac{1}{2} \underline{w}'B^{-1}\underline{w}$$

$$\Rightarrow B^{-1} = \alpha I + \beta \Phi' \Phi, \text{ ここで } \Phi = \begin{pmatrix} \phi_0(x_1) & \dots & \phi_{M-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_N) & \dots & \phi_{M-1}(x_N) \end{pmatrix}$$

②  $\varphi(w)$  の  $w$  に関する 1 次の項を比較すると  
 $\beta \sum_{n=1}^N t_n w' \phi(x_n) = \frac{1}{2} (a' B^{-1} w + w' B^{-1} a)$

$$\Rightarrow a = \beta B \Phi' t$$

レポート: 上の ①, ② を確かめよう.

結局 事後確率は

$$p(w|t, \alpha, \beta) = c' \exp \varphi(w) \\ = c' \exp \left\{ -\frac{1}{2} (w-a)' B^{-1} (w-a) \right\}$$

多次元ガウス分布 ( $c'$ : 規格化定数) !!

よ) 推定値は

$$w_{\text{MAP}} = a$$

Maximum a Posteriori (最大事後)

で与えられる

コメント: 最大推定と比較して, バイズ推定の方が over fitting が生じにくい.

## 時系列バイズ推定

$N$  のデータ  $(x_n, t_n)$   $n=1, \dots, N$  が順次手に入るとき  
 その都度  $w$  を学習 (推定) していく方法

アイデア: 事後確率を次のステップの事前確率に使う

$$p(w) \dots \text{事前 } p$$

$$p(t_1|w) \longrightarrow p(w|t_1) = c_1 p(t_1|w) p(w) \dots \text{ステップ } 1$$

( $x_1, t_1$ ) の尤度関数      事後      事前確率として使う

$$p(t_2|w) \longrightarrow p(w|t_1, t_2) = c_2 p(t_2|w) p(w|t_1) \dots \text{ステップ } 2$$

( $x_2, t_2$ ) の尤度関数

$$\vdots$$

←  $k-1$  ステップ後で得られる事後確率

$$p(t_k|w) \longrightarrow p(w|t_1, \dots, t_k) = c_k p(t_k|w) p(w|t_1, \dots, t_{k-1}) \dots \text{ステップ } k$$

## 今日の話題: 分類問題

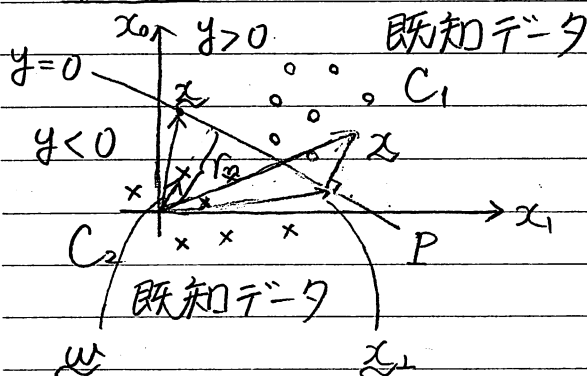
(Classification)

「入力  $x \in \mathbb{R}^D$  に対して  $K$  個の異なるクラス $C_k$ ,  $k=1, \dots, K$  のどれか1つを割り当てたい」記号:  $t(x) = (t_1, \dots, t_K)$  対  $x$  は  $C_j$  の  
 $t_j = 1, t_k = 0$  ( $k \neq j$ )  $\Leftrightarrow$  クラス  $k$  属する。問題設定: 既知の分類分けデータ  $(x_i, t_i)$ ,  $i=1, \dots, N$   
を用いて未知入力データの分類分けを行う方針: 既知データを用いて  $\mathbb{R}^D$  を  $K$  個の部分集合 $C_k$ ,  $k=1, \dots, K$  (同じ記号を使う) に分割

$$\mathbb{R}^D = \bigcup_{k=1}^K C_k$$

できればよい

## 非確率的方法

①  $K=2$ . $f(x) = w'x + w_0 = 0 \Leftrightarrow (D-1)$ 次元超平面を定める。  
 $P$ とする。 $D=2$ .既知データ  $x$  から  
 $(w, w_0)$  を推定 (学習)  
できればよい未知データ  $x_k$  に対しては.

$$f(x) \geq 0 \Rightarrow x \in C_1$$

$$f(x) < 0 \Rightarrow x \in C_2$$

とする。

## 幾何的考察

$$x_A, x_B \in P \Rightarrow f(x_A) = f(x_B)$$

$$\Rightarrow w'(x_A - x_B) = 0$$

よて  $\omega$  は  $P$  に直交するベクトル

・原点から  $P$  までの最短距離  $r_n$  は  $x \in P$  をとると

$$r_n = \frac{\omega'x}{\|\omega\|} = \frac{-\omega_0}{\|\omega\|}$$

で与えられる。

・  $x \in \mathbb{R}^D$  に対して  $x_{\perp} \in P$  を  $x$  の  $P$  上への射影とすると

$$x = x_{\perp} + r \frac{\omega}{\|\omega\|}, \quad r \in \mathbb{R} \text{ とかける.}$$

$$y(x) = \omega'x_{\perp} + r \frac{\omega'\omega}{\|\omega\|} + \omega_0 = r\|\omega\| \quad \therefore r = \frac{y(x)}{\|\omega\|}$$

①  $K > 2$  の場合.

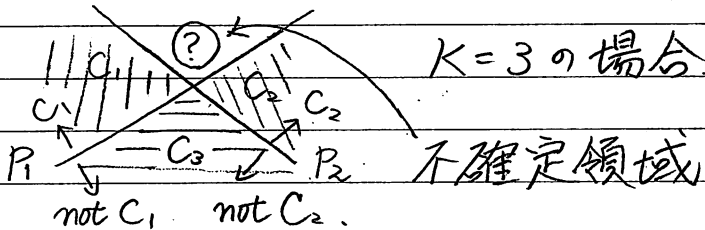
悪い分類の仕方 1

$K-1$  個の超平面  $P_k, k=1, \dots, K-1$  ( $\Leftrightarrow y_k(x) = 0$ )

を導入し

★  $y_k(x) \geq 0 \Leftrightarrow x \in C_k, k=1, \dots, K-1$

とする。



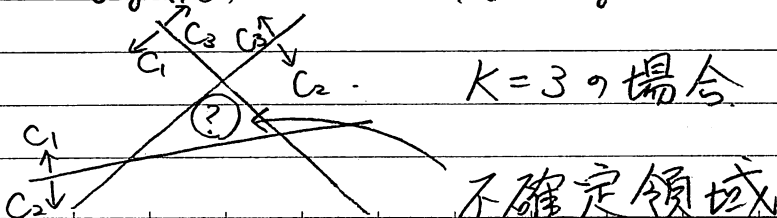
悪い分類の仕方 2.

$C_i, C_j$  の分類を行う超平面  $P_{ij} (i \neq j)$  を

$K C_2 = \binom{K}{2} = \frac{K(K-1)}{2}$  を導入

★  $y_{ij}(x) \geq 0 \Leftrightarrow x \in C_i$

$y_{ij}(x) < 0 \Leftrightarrow x \in C_j$



良い分類の仕方.

$K$ 本の線形関数を導入:

$$y_k(x) = \omega'_k x + \omega_{k0}, \quad k=1, \dots, K$$

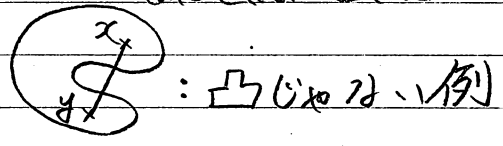
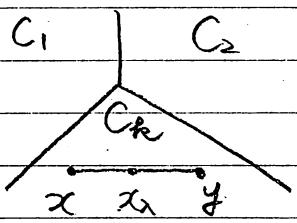
$x \in \mathbb{R}^D$   $K$  対して.

★  $y_k(x) \geq y_j(x)$  (おのれの  $k \neq j$ )  $\Leftrightarrow x \in C_k$   
で分類分けを考えろ.

補題

各部分集合  $C_k \subset \mathbb{R}^D$  は凸集合

(任意の2点  $x, y \in C_k$  対して  
それらを結ぶ直線も  $C_k$  に入っている)



(証明)  $x_A, x_B \in C_k$  とする. このとき  $x_A$  と  $x_B$  を  
結ぶ直線は.

$$\{x_\lambda = \lambda x_A + (1-\lambda)x_B \mid \lambda \in [0, 1]\}$$

と表される. 仮定から

$$y_k(x_A) \geq y_j(x_A) \quad (\forall j \neq k)$$
$$y_k(x_B) \geq y_j(x_B)$$

が成り立つ.

示したい事:  $\forall \lambda \in [0, 1]$

$$x_\lambda \in C_k \quad (\Leftrightarrow y_k(x_\lambda) \geq y_j(x_\lambda) \quad (\forall j \neq k))$$

$$y_k(x_\lambda) = \omega'_k (\lambda x_A + (1-\lambda)x_B) + \omega_{k0}$$
$$= \lambda \omega'_k x_A + (1-\lambda) \omega'_k x_B + \lambda \omega_{k0} + (1-\lambda) \omega_{k0}$$
$$= \lambda y_k(x_A) + (1-\lambda) y_k(x_B)$$
$$\geq \lambda y_j(x_A) + (1-\lambda) y_j(x_B)$$
$$= \lambda \omega'_j x_A + (1-\lambda) \omega'_j x_B + \omega_{j0} = y_j(x_\lambda) \quad \parallel$$



よって既知データ  $(x_i, t_i) \quad i=1, \dots, N$  を用いて

$$\tilde{W} = \begin{pmatrix} w_{10} & w_{20} & \dots & w_{k0} \\ \tilde{w}_1 & \tilde{w}_2 & \dots & \tilde{w}_k \end{pmatrix} = (\tilde{w}_1, \dots, \tilde{w}_k) \}_{D+1 \times}$$

を学習する.  $\tilde{w}_k' = (w_{k0}, w_{k1}, \dots, w_{kD})$

$\tilde{x}' = (1, x_1, \dots, x_D)$  とすると

$y_k(x) = \tilde{w}_k' \tilde{x}$  と書ける.

$(x_i, t_i) \quad i=1, \dots, N$   $k$  対して

$$\tilde{X} = \begin{pmatrix} \tilde{x}'_1 \\ \vdots \\ \tilde{x}'_N \end{pmatrix}, \quad T = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$$

$$\tilde{X}\tilde{W} = \begin{pmatrix} \tilde{x}'_1 \\ \vdots \\ \tilde{x}'_N \end{pmatrix} (\tilde{w}_1, \dots, \tilde{w}_k) = \begin{pmatrix} y_1(x_1) & \dots & y_k(x_1) \\ \vdots & \ddots & \vdots \\ y_1(x_N) & \dots & y_k(x_N) \end{pmatrix}$$

$$\tilde{X}\tilde{W} - T = \begin{pmatrix} y_1(x_1) - t_{11} & \dots & y_k(x_1) - t_{1k} \\ \vdots & \ddots & \vdots \\ y_1(x_N) - t_{N1} & \dots & y_k(x_N) - t_{Nk} \end{pmatrix}$$

### 補題

$(y_1, \dots, y_k)$   $k$  対して

$$y_k > y_i \quad (\forall i \neq k) \iff \underbrace{\sum_{j \neq k} y_j^2 + (y_k - 1)^2}_{A_k} < \underbrace{\sum_{j \neq i} y_j^2 + (y_i - 1)^2}_{A_i} \quad (\forall i \neq k)$$

$$\text{(証)} \quad x > y \iff (x-1)^2 + y^2 < x^2 + (y-1)^2 \quad \textcircled{1}$$

$$A_i - A_k = y_k^2 + (y_i - 1)^2 - y_i^2 - (y_k - 1)^2$$

$$\textcircled{1} \text{より} \quad A_i - A_k > 0 \iff y_k > y_i \quad //$$

任意の  $i$  行目  $k$  に対して

$\star i \quad \sum_{j=1}^k (y_j(x_i) - t_{ij})^2$  が最小

$k$  なるよりの  $\tilde{W}$  の既知データの正しい分類を与える

エラー関数:

$$E(\tilde{W}) = \frac{1}{2} \text{Tr} \{ (X\tilde{W} - T)(X\tilde{W} - T)'\}$$

$$= \frac{1}{2} \sum_{i=1}^N \star_i$$

の最小点,  $\tilde{W}$  を求める.

Remark:  $E(\tilde{W})$  は  $\tilde{W}$  に関して 2 次

$$\frac{\partial}{\partial \tilde{W}} E(\tilde{W}) = 0 \text{ なる } \tilde{W} ?$$

$A, B$  を  $n \times n$  行列

$$\left. \begin{aligned} \frac{\partial}{\partial A} \text{Tr}(A'BA) &= (B+B')A \\ \frac{\partial}{\partial A} \text{Tr}(A'B) &= B \end{aligned} \right\} (\$)$$

(\\$) を使くと

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

よってこの  $\tilde{W}_*$  を用いることで  $K$  個の分類

$C_k$   $k=1, \dots, K$  を定める線形関数

$$\begin{pmatrix} y_1^*(x) \\ \vdots \\ y_k^*(x) \end{pmatrix} = \tilde{W}_*' x$$

が決定できた. 未知データ  $x \in \mathbb{R}^D$  に対しては

$$y_k^*(x) \geq y_j^*(x), \text{ 全ての } j \neq k$$

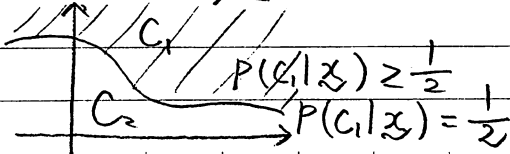
なら  $k$  を採って  $x \in C_k$  とすればよい.

今日の話題: 分類問題 (確率的手法)

話を簡単にする:  $K=2$  (2種類のクラス  $C_1, C_2$ )

要は...

$P(C_1|x)$  を求めて  $P(C_1|x) = 1/2$  の曲面を  $\mathbb{R}^d$  内で定める. これより  $x \in C_1 \Leftrightarrow P(C_1|x) \geq 1/2$



手順:

 $(k=1,2)$ 

$$\textcircled{1} p(x|C_k) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)\right\}$$

$$= N(x|\mu_k, \Sigma) \text{ 多次元ガウス分布}$$

$$p(C_1) = \pi, p(C_2) = 1 - \pi \quad 0 \leq \pi \leq 1$$

と仮定し、既知データ  $(x_i, t_i) \quad i=1, \dots, N$  を用いて  
パラメータ  $\mu_1, \mu_2, \Sigma, \pi$  を学習

$$\text{ここで} \quad t_i = \begin{cases} 1 & x_i \in C_1 \\ 0 & x_i \in C_2 \end{cases}$$

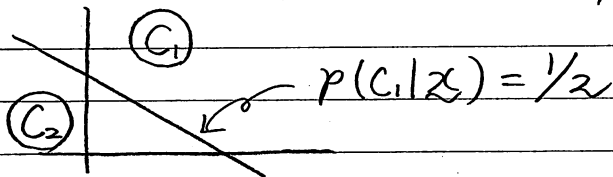
② バイズの法則を用いて

$$p(C_1|x) = \frac{p(x|C_1)p(C_1)}{p(x)}$$

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

を求める。(当然、 $p(C_2|x) = 1 - p(C_1|x)$ )

③  $p(C_1|x) = 1/2$  を求める超曲面 (実は超平面  
 $k$  なら) が  $C_1$  と  $C_2$  の境界として得られる。



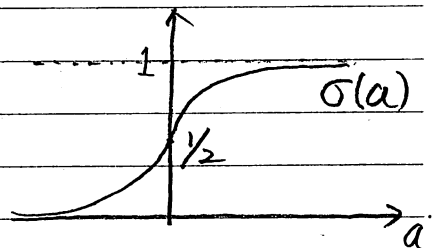
②, ③ を先に説明 (パラメータが推定できたとする)

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

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

$$\text{ここで} \quad a = \log \frac{p(x|C_1)p(C_1)}{p(x|C_2)p(C_2)}$$

$\sigma(a)$ : logistic sigmoid 関数



ここで  $P(x|C_k)$ ,  $P(C_k)$  を代入すると

$$\begin{aligned} a &= \log p(x|C_1) - \log p(x|C_2) + \log p(C_1)/p(C_2) \\ &= -\frac{1}{2}(x-\mu_1)' \Sigma^{-1}(x-\mu_1) + \frac{1}{2}(x-\mu_2)' \Sigma^{-1}(x-\mu_2) \\ &\quad + \log \frac{\pi}{1-\pi}. \end{aligned}$$

$x$  の 2 次の項は打ち消される。

$$\begin{aligned} &= x' \Sigma^{-1} \mu_1 - x' \Sigma^{-1} \mu_2 - \frac{1}{2} \mu_1' \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2' \Sigma^{-1} \mu_2 \\ &\quad + \log \frac{\pi}{1-\pi} \\ &= \underline{w}' x + w_0 \end{aligned}$$

$$\text{ここで } w = \Sigma^{-1}(\mu_1 - \mu_2)$$

$$w_0 = -\frac{1}{2} \mu_1' \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2' \Sigma^{-1} \mu_2 + \log \frac{\pi}{1-\pi}$$

$$\text{よって } p(C_1|x) = \sigma(w'x + w_0)$$

手順③から

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

つまり  $w'x + w_0 = 0$  : 超平面 !!

が分類分けを定める境界

あとは手順①のパラメータ推定を考えればよい

既知データ:  $\{x_n, t_n\} \quad n=1, \dots, N \quad t_n = \begin{cases} 1 & (x_n \in C_1) \\ 0 & (x_n \in C_2) \end{cases}$

$$\begin{aligned} x_n \in C_1 &\Rightarrow p(x_n, t_n=1) = p(C_1) p(x|C_1) \\ &= \pi N(x_n | \mu_1, \Sigma) \end{aligned}$$

$$\begin{aligned} x_n \in C_2 &\Rightarrow p(x_n, t_n=0) = p(C_2) p(x|C_2) \\ &= (1-\pi) N(x_n | \mu_2, \Sigma) \end{aligned}$$

よって  $t = (t_1, \dots, t_N)$ ,  $X = (x_1, \dots, x_N)$  とすると尤度関数は

$$p(t, X | \mu_1, \mu_2, \Sigma)$$

$$= \pi_{n=1}^N [\pi N(x_n | \mu_1, \Sigma)]^{t_n} [(1-\pi) N(x_n | \mu_2, \Sigma)]^{1-t_n}$$

とみる。

最大点  $(\pi^*, \mu_1^*, \mu_2^*, \Sigma^*)$  を求める。(最大推定)

$$\log p(t, X | \mu_1, \mu_2, \Sigma)$$

$$= \sum_{n=1}^N [t_n \{ \log \pi + \log N(x_n | \mu_1, \Sigma_1) \} + (1-t_n) \{ \log(1-\pi) + \log N(x_n | \mu_2, \Sigma_2) \}]$$

$$\frac{\partial}{\partial \pi} \log p(x | \star) = \sum_{n=1}^N \left( \frac{t_n}{\pi} - \frac{1-t_n}{1-\pi} \right) = 0$$

$$(1-\pi) \sum_{n=1}^N t_n = \pi \sum_{n=1}^N (1-t_n)$$

$$\sum_{n=1}^N t_n = n\pi \Rightarrow \pi^* = \frac{1}{N} \sum_{n=1}^N t_n$$

$$\frac{\partial}{\partial \mu_1} \log p(x | \star) = \frac{\partial}{\partial \mu_1} \sum_{n=1}^N t_n \log N(x_n | \mu_1, \Sigma_1)$$

$$= \frac{\partial}{\partial \mu_1} \sum_{n=1}^N t_n \left( -\frac{1}{2} (x_n - \mu_1)' \Sigma_1^{-1} (x_n - \mu_1) + \log(\mu_1 \text{以外}) \right)$$

$$= \sum_{n=1}^N t_n \left[ \frac{\partial}{\partial \mu_1} (\mu_1' \Sigma_1^{-1} x_n) - \frac{\partial}{\partial \mu_1} \left( \frac{1}{2} \mu_1' \Sigma_1^{-1} \mu_1 \right) \right]$$

$$= \sum_{n=1}^N t_n \Sigma_1^{-1} x_n - \sum_{n=1}^N t_n \Sigma_1^{-1} \mu_1 = 0$$

$$\sum_{n=1}^N t_n x_n = \left( \sum_{n=1}^N t_n \right) \mu_1$$

$$\mu_1^* = \frac{1}{\sum_{n=1}^N t_n} \sum_{n=1}^N t_n x_n$$

同様く

$$\mu_2^* = \frac{1}{\sum_{n=1}^N (1-t_n)} \sum_{n=1}^N (1-t_n) x_n$$

最後く  $\Sigma_1^*$  については  $\frac{\partial}{\partial \Sigma_1} \log p(x | \star) = 0$  より

$$\Sigma_1^* = \frac{1}{N_1} S_1 + \frac{1}{N_2} S_2$$

ここで  $S_1, S_2$  (行列) は

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (x_n - \mu_1^*) (x_n - \mu_1^*)'$$

$$S_2 = \frac{1}{N_2} \sum_{n \in C_2} (x_n - \mu_2^*) (x_n - \mu_2^*)'$$

$$N_1 = \sum_{n=1}^N t_n, \quad N_2 = \sum_{n=1}^N (1-t_n)$$

# Support Vector Machine (SVM) 入門

## 復習 (2分類問題)

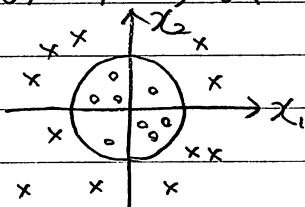
既知データを用いて超平面

$$y(x) = w'x + w_0 = 0 \quad x \in \mathbb{R}^d$$

を学習させて (つまり)  $w, w_0$  を推定)

未知データの分類を行う

既知データが



⇒ 打用策

$$y(x) = x_1^2 + x_2^2 - r^2$$

$$= w' \phi(x) + w_0$$

$$w' = (1, 1), \phi(x) = (x_1^2, x_2^2)$$

$$w_0 = -r^2$$

とすることで既知データの分類は可能になる。

の場合には平面での分類は不可能

## 今後一般化

$$(1) y(x) = w' \phi(x) + w_0$$

$\phi(x) = (\phi_1(x), \dots, \phi_n(x))$  はある非線形変換

の形で分類問題を考えよう (つまり曲面による

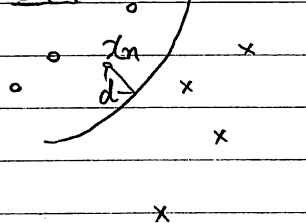
$\mathbb{R}^D$  の分割)

既知データを

$$t_n = \begin{cases} 1 & y(x_n) > 0 \Leftrightarrow C_1 \text{ の領域} \\ -1 & y(x_n) < 0 \Leftrightarrow C_2 \text{ " } \end{cases}$$

としよう  $|y(x_n)| = t_n y(x_n)$  に注意

$\mathbb{R}^D$   $y(x) = 0$



0 と x を分類する曲面  $y(x) = 0$  は連続的に存在するが、今後の未知データを最も正確に分類できようのは、

★「最も近い既知データとの距離が最大になる曲面」

$$d = \frac{|y(x_n)|}{\|\underline{w}\|} = \frac{t_n y(x_n)}{\|\underline{w}\|} = \frac{t_n (\underline{w}' \phi(x_n) + w_0)}{\|\underline{w}\|} \quad \text{よ') ☆ の}$$

$$(\$) \operatorname{argmax}_{\underline{w}, w_0} \left\{ \frac{1}{\|\underline{w}\|} \min_n [t_n (\underline{w}' \phi(x_n) + w_0)] \right\}$$

で与えられる。

{\*} の  $\underline{w} \rightarrow \kappa \underline{w}$ ,  $w_0 \rightarrow \kappa w_0$  ( $0 \neq \kappa \in \mathbb{R}$ ) で不変なので  
最も曲面に近しい点  $x_n$  で

$$t_n (\underline{w}' \phi(x_n) + w_0) = 1$$

と規格化しておく。すると

$$t_i (\underline{w}' \phi(x_n) + w_0) \geq 1 \quad i=1, \dots, N$$

これよ') (\$) の

$$\operatorname{argmax}_{\underline{w}, w_0} \left\{ \frac{1}{\|\underline{w}\|} \min_n \left[ \frac{1}{1} \right] \right\} = \operatorname{argmax}_{\underline{w}, w_0} \frac{1}{\|\underline{w}\|},$$

つまり

$$\operatorname{argmin}_{\underline{w}, w_0} \frac{1}{2} \|\underline{w}\|^2$$

の解として曲面が求まる。

まとめると 既知データ  $\{x_n, t_n\}$  の曲面 (1) による分類の

制約条件:  $t_i (\underline{w}' \phi(x_i) + w_0) \geq 1 \quad i=1, \dots, N$

の下で

$$\operatorname{argmin}_{\underline{w}, w_0} \frac{1}{2} \|\underline{w}\|^2$$

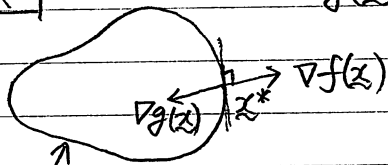
を求めれば良い。

制約条件付き極値問題  $\Leftarrow$  ラグランジュの未定乗数法

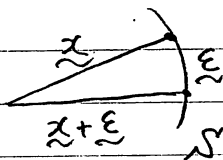
復習 ラグランジュ未定乗数法

① 制約条件:  $g(x) = 0 \quad x \in \mathbb{R}^D$   
 の下で関数  $f(x) \in \mathbb{R}^D$  を最大にする  
 点  $x^*$  を求めたい。 } Q1

$\mathbb{R}^D$   $g(x) = 0 \leftarrow \mathbb{R}^D$  の  $(D-1)$  次元曲面  
を定める。  $\downarrow$   
 $S$  とする



$$S = \{g(x) = 0\}$$



補題:  $x \in S$  に対して

$$\nabla g(x) = \left( \frac{\partial g}{\partial x_1}, \dots, \frac{\partial g}{\partial x_D} \right) \leftarrow g \text{ が最も} \\ \text{増える方向}$$

は  $S$  に直交する

証)  $x + \varepsilon$  を  $S$  上の  $x$  に十分近い点とするとき  
(テイラー展開より)

$$g(x + \varepsilon) = g(x) + \varepsilon' \nabla g(x) + O(\varepsilon^2)$$

$$0 = 0$$

$$\varepsilon' \nabla g(x) = 0$$

$\varepsilon$  は  $x$  での  $g(x) = 0$  の接ベクトルなので  
 $\nabla g(x)$  は  $S$  に直交する。

一方  $f(x)$  が  $x^* \in S$  で極大値をとるならば  
 $\nabla f(x^*)$  は  $S$  に直交

よって

$$\nabla f(x^*) + \lambda \nabla g(x^*) = 0$$

とすれば  $0 \neq \lambda \in \mathbb{R}$  が存在する ( $\because \nabla f$  と  $\nabla g$  は平行)

$$L(x, \lambda) = f(x) + \lambda g(x) : \text{ラグランジュ関数}$$

ラグランジュの未定乗数

$Q$  が極大値を持つ必要条件

$$\nabla_x L = 0, \quad \frac{\partial}{\partial \lambda} L = 0$$

これを満たす  $(x^*, \lambda^*)$  を求める方法

ラグランジュ未定乗数法



② 制約条件:  $g(x) \geq 0, x \in \mathbb{R}^D$   
 の下で  $f(x) \in \mathbb{R}$  を極大にする点  $x^*$   
 を求める } Q2

ラグランジュ関数

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

$\nabla f(x^*)$  (Bと等しい)  
 $\nabla g(x) \leftarrow g(x) < 0$   
 $g(x) > 0$   
 $g(x) = 0$

場合分け

(A)  $g(x^*) > 0$

この場合 制約条件の無意味

$$\Rightarrow \lambda = 0, \frac{\partial L}{\partial x} = 0$$

(B)  $g(x^*) = 0$ , Q1と同じ

$$\text{よって } 0 \neq \lambda^* \in \mathbb{R} \text{ に対して } \nabla_x L(x^*, \lambda^*) = \frac{\partial L}{\partial x}(x^*, \lambda^*) = 0$$

$$\text{ただし } \nabla f(x^*) = -\lambda^* \nabla g(x^*) \text{ より } \lambda^* > 0 \text{ でありければ}$$

なりなさい。

よってまとめると: ← 場合分けに対応

$$g(x) \geq 0, \lambda \geq 0, \lambda g(x) = 0 \text{ の下で}$$

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

の極値を与えら  $(x^*, \lambda^*)$  を求める方法が Q2

に対応するラグランジュ未定乗数法

複数個の制約条件  $g_j(x) = 0, j = 1, \dots, J$

$h_k(x) \geq 0, k = 1, \dots, K$  の場合

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

$$\text{(ただし } \mu_k \geq 0, \mu_k h_k(x) = 0 \text{)}$$

とすればよい。

解こう

$$\textcircled{1} \begin{cases} L(w, w_0, a) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N a_n \{ \tanh(w' \phi(z_n) + w_0) - 1 \} \\ \text{として } a_n \geq 0, \tanh y(z_n) \geq 1, a_n (\tanh y(z_n) - 1) = 0 \end{cases}$$

の下で  $L(w, w_0, a)$  の極値を探す

$$\frac{\partial L}{\partial w} = w - \sum_{n=1}^N a_n t_n \phi(x_n) = 0 \Rightarrow w^* = \sum_{n=1}^N a_n t_n \phi(x_n)$$

$$\frac{\partial L}{\partial w_0} = -\sum_{n=1}^N a_n t_n = 0$$

$L(w, w_0, a)$  を代入すると

$$\Rightarrow \tilde{L}(a) = \sum_{n=1}^N a_n + \frac{1}{2} \left( \sum_{n=1}^N a_n t_n \phi(x_n)' \right) \left( \sum_{m=1}^N a_m t_m \phi(x_m) \right) - \sum_{n=1}^N a_n t_n \phi(x_n)$$

$$\textcircled{2} = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(x_n, x_m)$$

ここで  $k(x, x') = \phi(x)' \phi(x)$  : カーネル

$$a_n \geq 0, \sum_{n=1}^N a_n t_n = 0$$

の下で極値  $\tilde{L}(a^*)$  が求まれば  $w^* = \sum_{n=1}^N a_n^* t_n \phi(x_n)$  として完了

→ 二次計画問題

(非線形)最適化問題

多くの数値解法が知られている。(MATLAB の関数もある)

⇒  $a^*$  を数値的に求めらることは可能!

学習解  $w^*$  の性質

分類を定める曲面  $y=0$  は

$$y(x) = w^* \phi(x) + w_0$$

$$= \sum_{n=1}^N a_n t_n \phi(x_n)' \phi(x) + w_0$$

$$= \sum_{n=1}^N a_n t_n k(x_n, x) + w_0 = 0$$

と与えられる

ここで全ての既知データに対して

$$a_n = 0 \quad \text{or} \quad t_n y(x_n) = 1$$

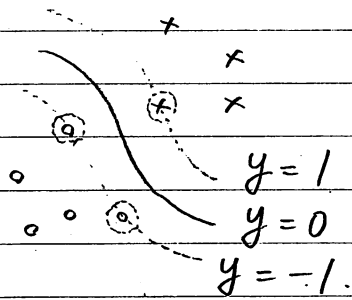
が満たされている。 (:: ①  $a_n \geq 0, a_n(t_n y(x_n) - 1) = 0$  )

よって

$$y(x) = \sum_{n \in S} a_n t_n k(x_n, x) + w_0$$

$$S := \{n \mid t_n y(x_n) = 1\}$$

分類曲面  $y(x) = 0$  に最も近くなる既知データ



⇒ 未知データの分類の計算量が大幅に減少

(左の例だと

3つの○でOK)

○ : サポート

まとめると

サポートベクトルマシンの特徴

「 $n \in \mathcal{S}$  は  $n$  (サポート) のみを用いて未知データの分類を行う」

# 資料 1 (参考書 p7~8)

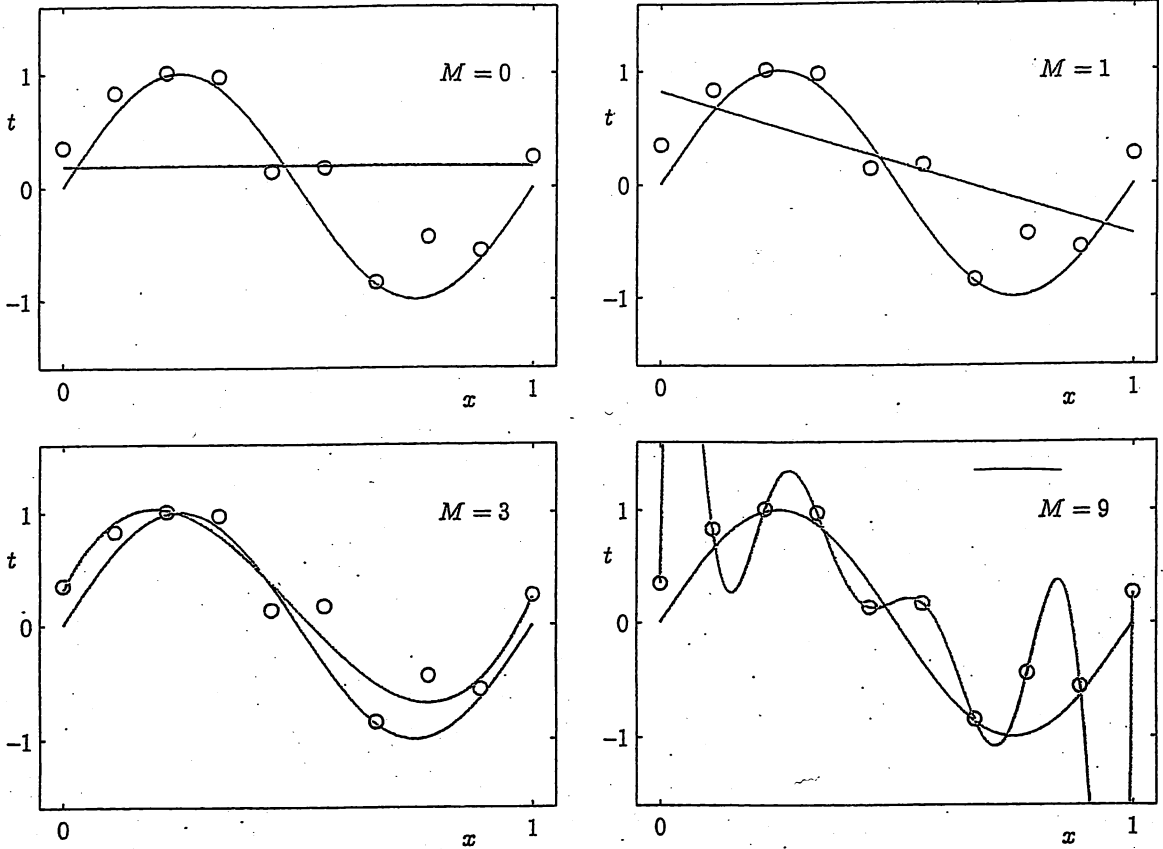


Figure 1.4 Plots of polynomials having various orders  $M$ , shown as red curves, fitted to the data set shown in Figure 1.2.

**Table 1.1** Table of the coefficients  $w^*$  for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

|         | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$     |
|---------|---------|---------|---------|-------------|
| $w_0^*$ | 0.19    | 0.82    | 0.31    | 0.35        |
| $w_1^*$ |         | -1.27   | 7.99    | 232.37      |
| $w_2^*$ |         |         | -25.43  | -5321.83    |
| $w_3^*$ |         |         | 17.37   | 48568.31    |
| $w_4^*$ |         |         |         | -231639.30  |
| $w_5^*$ |         |         |         | 640042.26   |
| $w_6^*$ |         |         |         | -1061800.52 |
| $w_7^*$ |         |         |         | 1042400.18  |
| $w_8^*$ |         |         |         | -557682.99  |
| $w_9^*$ |         |         |         | 125201.43   |

