

3. 線形回帰および識別

線形回帰のモデル

正則化項の導入

L2正則化

L1正則化

正則化項のBayes的解釈

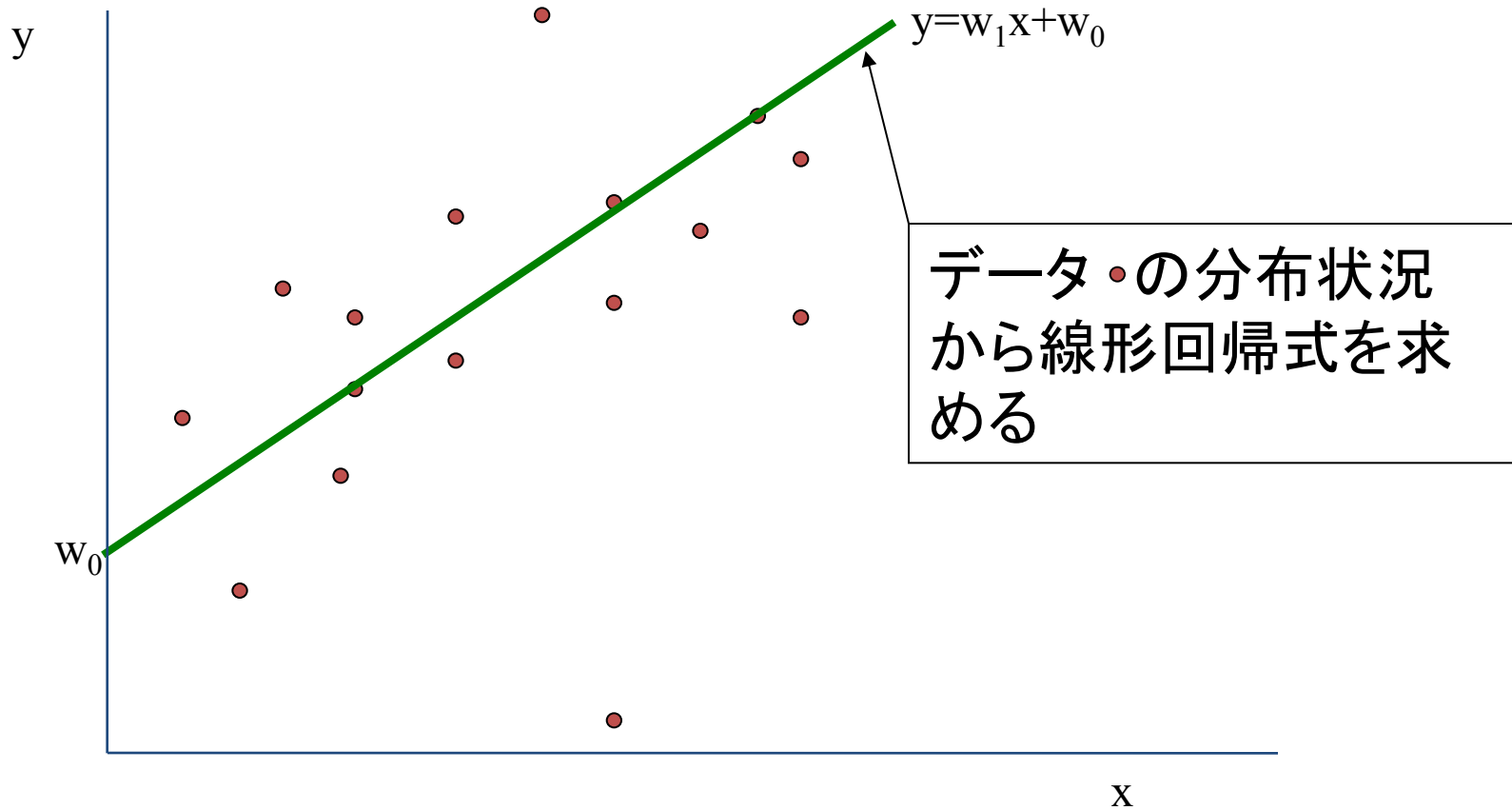
線形識別

生成モデルを利用した識別

2乗誤差最小化の線形識別の問題点

by 中川裕志(東京大学)

線形モデル



線形モデル

- 入力ベクトル: \mathbf{x} から出力: y を得る関数が \mathbf{x} の線形関数 (\mathbf{w} と \mathbf{x} の内積)

$$y = \langle \mathbf{x}, \mathbf{w} \rangle = \sum_{i=0}^K w_i x_i \quad \text{ただし、} \mathbf{x} = [1, x_1, \dots, x_K]^T, \mathbf{w} = [w_0, w_1, \dots, w_K]^T$$

- 一般に観測データはノイズを含んでいる。つまり

$$y = \langle \mathbf{x}, \mathbf{w} \rangle + \varepsilon \quad \varepsilon \text{ はノイズで } N(0, \sigma^2) \text{ と考える。}$$

- 得られた N 個の観測データの組 (\mathbf{y}, \mathbf{X}) に対して最適な \mathbf{w} を推定する。
- そこで、 \mathbf{y} と $\mathbf{X}\mathbf{w}$ の2乗誤差を最小化するように \mathbf{w} を選ぶ。

2乗誤差の最小化

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1K} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & \cdots & x_{NK} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{bmatrix}$$

$$\mathbf{w} \text{ の推定値 } \hat{\mathbf{w}} = \arg \min_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\frac{\partial (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = 0 \quad \text{を解く と } \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

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

➤ 正規方程式 と呼ばれる基本式

補遺: 正規方程式の導出

$$(y - Xw)^T (y - Xw) = (y^T - w^T X^T)(y - Xw) = y^T y - w^T X^T y - y^T Xw + w^T X^T Xw$$

$$\frac{\partial (y - Xw)^T (y - Xw)}{\partial w} = -\frac{\partial w^T X^T y}{\partial w} - \frac{\partial y^T Xw}{\partial w} + \frac{\partial w^T X^T Xw}{\partial w} = 0 \quad (1)$$

$$\frac{\partial x^T a}{\partial x} = a \text{ より } \frac{\partial w^T X^T y}{\partial w} = X^T y \quad \frac{\partial a^T x}{\partial x} = a \text{ より } \frac{\partial y^T Xw}{\partial w} = (y^T X)^T = X^T y$$

$$\frac{\partial w^T X^T Xw}{\partial w} = \frac{\partial w^T (X^T Xw)}{\partial w} + \frac{\partial (w^T X^T X)w}{\partial w} = X^T Xw + (w^T X^T X)^T = 2X^T Xw$$

$$\Rightarrow (1) = -2X^T (y - Xw) = -2X^T y + 2X^T Xw = 0 \quad \Rightarrow X^T y = X^T Xw$$

$$\Rightarrow w = (X^T X)^{-1} X^T y$$

cf 行列で微分する場合のchain rule $\frac{\partial f(g(x))}{\partial x} = \frac{\partial g(x)}{\partial x} \frac{\partial f(g(x))}{\partial g(x)}$ を使えば

$$\frac{\partial (y - Xw)^T (y - Xw)}{\partial w} = \frac{\partial (y..)}{\partial w} \frac{\partial (y..) ^T (y..)}{\partial (y..)} + \frac{\partial (y..) ^T}{\partial w} \frac{\partial (y..) ^T (y..)}{\partial (y..) ^T}$$

$$= -X^T (y - Xw) - X^T (y - Xw) = -2X^T (y - Xw)$$

正規方程式を解く簡単な例

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \text{正規方程式} \quad \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \text{は}$$

$$\begin{bmatrix} 1 & \cdots & 1 \\ x_1 & \cdots & x_N \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} 1 & \cdots & 1 \\ x_1 & \cdots & x_N \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N x_i y_i \end{bmatrix} \quad \Leftrightarrow \quad \mathbf{X}^T \mathbf{X} \quad \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2} \begin{bmatrix} \sum_{i=1}^N x_i^2 & -\sum_{i=1}^N x_i \\ -\sum_{i=1}^N x_i & N \end{bmatrix} \Rightarrow w_1 = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2}$$

$$w_0 = \frac{\sum_{i=1}^N x_i^2 \sum_{i=1}^N y_i - \sum_{i=1}^N x_i \sum_{i=1}^N x_i y_i}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2} = \frac{\sum_{i=1}^N y_i}{N} - \frac{\sum_{i=1}^N x_i}{N} \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2} = w_0 = \frac{1}{N} \sum_{i=1}^N y_i - \frac{1}{N} w_1 \sum_{i=1}^N x_i$$

用語：誤差、損失、目的関数

- 線形モデルで最小化したかったのは2乗誤差
- 真のモデルにおける値(2乗誤差における y)と予測値(2乗誤差における Xw)の差異を表す関数を損失関数(単に損失)あるいは $Loss$ と呼び、 L で表すことが多い。
- 上記のような最適化問題において最小化(一般的には最適化)したい関数を目的関数と呼ぶ。
- 線形モデルの2乗誤差最小化では
2乗誤差＝損失＝目的関数

線形モデルの一般化

$$y = \langle \varphi(\mathbf{x}), \mathbf{w} \rangle \quad \varphi(\mathbf{x}) = [1, \phi_1(\mathbf{x}), \dots, \phi_K(\mathbf{x})]^T$$

基底関数 重み

N個の観測データ (y, \mathbf{x}) に対して $y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \quad \varphi(\mathbf{x}) = \begin{pmatrix} \phi_1(\mathbf{x})^T \\ \vdots \\ \phi_N(\mathbf{x})^T \end{pmatrix}$

$(y, \varphi(\mathbf{x}))$ が得られたとすると、2乗誤差を最小化する \mathbf{w} は前を同じく以下の通りだが、少し別の見方で解く。

$$\hat{\mathbf{w}} = (\varphi(\mathbf{x})^T \varphi(\mathbf{x}))^{-1} \varphi(\mathbf{x})^T \mathbf{y}$$

基底関数の例

$$\phi_j(x) = x^j \quad : \text{polynomial}$$

$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\} \quad : \text{Gaussian}$$

$$\phi_j(x) = \frac{1}{1 + \exp(-(x - \mu_j)/s)} \quad : \text{sigmoidal}$$

$$\phi_j(x) = \exp\left(2\pi i \frac{xj}{m}\right) \quad (m : \text{even}) : \text{Fast Fourier}$$

正規方程式を求める別の方法

$$y = \langle \phi(\mathbf{x}), \mathbf{w} \rangle + \varepsilon \quad \varepsilon = N(0, \beta^{-1}) \quad \beta = \sigma^{-2} \text{を精度と呼ぶ.}$$

$$p(y | \mathbf{x}, \mathbf{w}, \beta) = N(y | \langle \phi(\mathbf{x}), \mathbf{w} \rangle, \beta^{-1})$$

- $\{\mathbf{x}(\text{ベクトル}), y\}$ が観測データ(training data)
- \mathbf{w}, β を決定する、即ち $p(y | \mathbf{x}, \mathbf{w}, \beta)$ を最大化)
- N 組の i.i.d. 観測データすなわち教師データがあるとする。

$$\mathbf{y} = (y_1, \dots, y_N)^T \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{bmatrix}$$

すると次のページのように $p(y | \mathbf{x}, \mathbf{w}, \beta)$ が書ける。

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N N(y_i \mid \langle \phi(\mathbf{x}_i), \mathbf{w} \rangle, \beta^{-1})$$

両辺のlogをとる

$$\log p(\mathbf{y} \mid \mathbf{w}, w_0, \mathbf{X}, \beta) = \frac{N}{2} \log \beta - \frac{N}{2} \log 2\pi - \beta L(\mathbf{w})$$

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \langle \phi(\mathbf{x}_i), \mathbf{w} \rangle)^2$$

$\log p(y|\mathbf{w}, \mathbf{X}, \beta)$ を \mathbf{w}, β について最大化したい。まず、 \mathbf{w} について最大化する。

$$\frac{\partial \log p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}, \beta)}{\partial \mathbf{w}} = \beta \sum_{i=1}^N \varphi(\mathbf{x}_i) (y_i - \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle) = 0$$

\Rightarrow

$$\sum_{i=1}^N \varphi(\mathbf{x}_i) y_i - \sum_{i=1}^N \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^T \mathbf{w} = 0$$

\Rightarrow

$$\boldsymbol{\varphi}(\mathbf{X})^T \cdot \mathbf{y} = (\boldsymbol{\varphi}(\mathbf{X})^T \boldsymbol{\varphi}(\mathbf{X})) \mathbf{w}$$

\Rightarrow

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

$$\boldsymbol{\varphi}(\mathbf{x}) = \begin{pmatrix} \varphi(\mathbf{x}_1)^T \\ \vdots \\ \varphi(\mathbf{x}_N)^T \end{pmatrix}$$

バイアス w_0 の部分だけに注目してみると

- 対数近似関数から最適な w_0 を によって求めると

$$\frac{\partial L(\mathbf{w})}{\partial w_0} = \frac{\partial \sum_{i=1}^N \left(y_i - (1, \phi_1(\mathbf{x}_i), \dots, \phi_K(\mathbf{x}_i)) \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{pmatrix} \right)^2}{\partial w_0} = \frac{\partial \sum_{i=1}^N \left(y_i - (\phi_1(\mathbf{x}_i), \dots, \phi_K(\mathbf{x}_i)) \begin{pmatrix} w_1 \\ \vdots \\ w_K \end{pmatrix} - w_0 \right)^2}{\partial w_0}$$
$$= -2 \sum_{i=1}^N \left(y_i - (1, \phi_1(\mathbf{x}_i), \dots, \phi_K(\mathbf{x}_i)) \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{pmatrix} \right) \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{pmatrix} = -2 \sum_{i=1}^N \left(y_i - (\phi_1(\mathbf{x}_i), \dots, \phi_K(\mathbf{x}_i)) \begin{pmatrix} w_1 \\ \vdots \\ w_K \end{pmatrix} - w_0 \right) = 0$$

\Rightarrow

$$w_0 = \frac{1}{N} \sum_{i=1}^N y_i - \frac{1}{N} \sum_{j=1}^K w_j \left(\sum_{i=1}^N \phi_j(\mathbf{x}_i) \right)$$

y の平均

基底関数の学習データの平均の w
重み付き和

精度 β を求める。

$\log p(y|\mathbf{w}, X, \beta)$ を β に対して最大化

ただし、 \mathbf{w} は最適化されたものを用いる

$$\frac{\partial \log p(\mathbf{y} | \hat{\mathbf{w}}, \mathbf{X}, \beta)}{\partial \beta} = \frac{N}{2\beta} - L(\hat{\mathbf{w}})$$

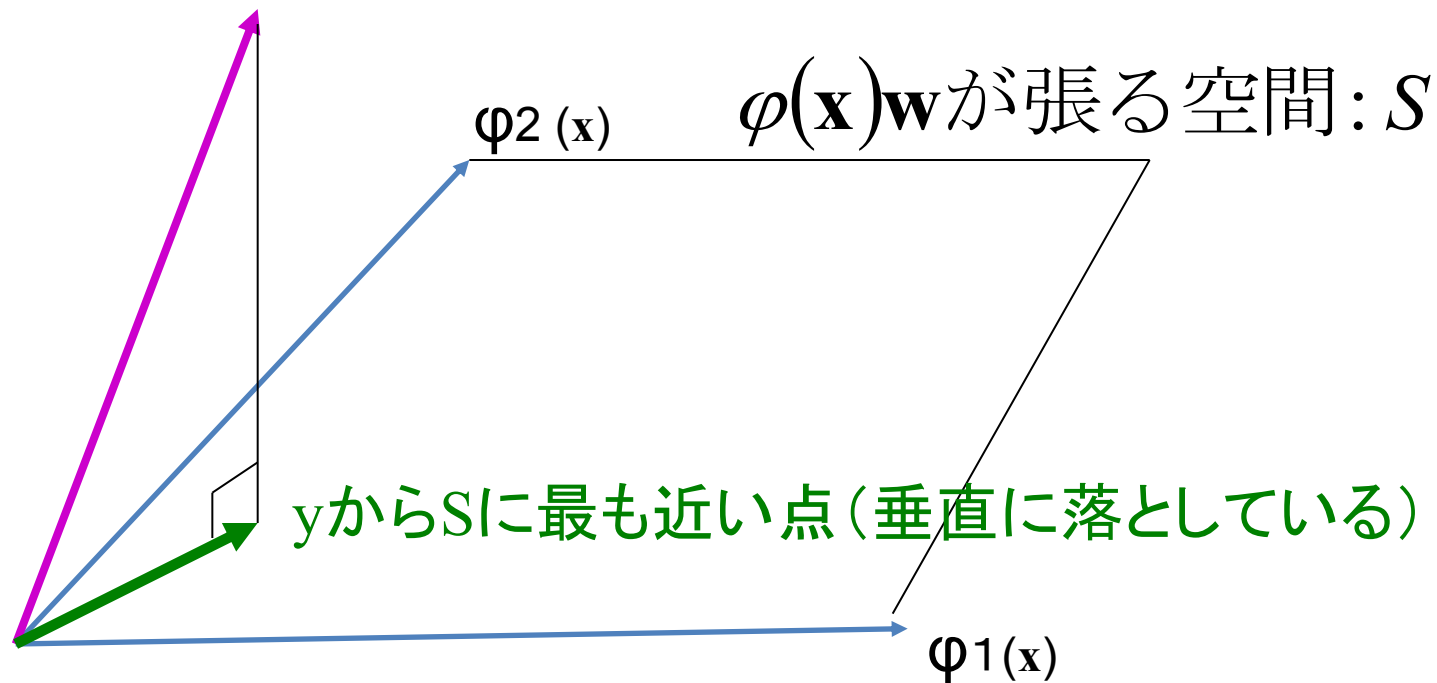
$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \boldsymbol{\varphi}(\mathbf{x}_i) \hat{\mathbf{w}})^2$$

$$\hat{\beta}^{-1} = \frac{1}{N} \sum_{i=1}^N (y_i - \boldsymbol{\varphi}(\mathbf{x}_i) \hat{\mathbf{w}})^2$$

y の予測値と観測された値の差の2乗の平均

幾何学的イメージ

新規データ: y



計算の効率化

- 大きなdata setsに対して $\hat{\mathbf{w}} = (\boldsymbol{\phi}^T \boldsymbol{\phi})^{-1} \boldsymbol{\phi}^T \mathbf{y}$ の右辺第1項の逆行列計算量が問題
- 特にデータの次元Nに対して $O(N^3)$ なので高次元だと大変
- 定石は、コレスキー分解 $O(N^2)$ して上/下半3角行列で表現される連立方程式を2回解く
- $L(\mathbf{w})$ を最小化するような \mathbf{w} の数値計算

$$\begin{aligned}\mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - \eta \nabla L(\mathbf{w}) \\ &= \mathbf{w}^{(\tau)} - \eta (y_n - \phi(\mathbf{x}_n) \mathbf{w}^{(\tau)}) \phi(\mathbf{x}_n)^T\end{aligned}$$

目的関数(すなわち損失 $L(\mathbf{w})$)の減る方向へ進む(— gradientを \mathbf{w} に加える)方法をgradient descent は呼ばれ、最適化における基本的数値計算法である。

正則化項の導入

- モデルを複雑にするほど学習データにはよく合致するが、学習データ以外のデータには弱いという過学習を起こす。
- 過学習を抑えるために、損失関数に正則化項を導入。
- 正則化項にはモデルをできるだけ単純化する方向に作用する。
 - データが高次元の場合には次元削減効果あり。

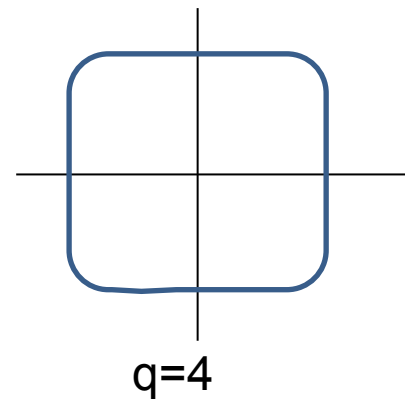
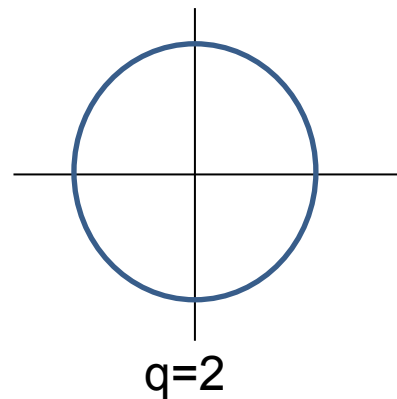
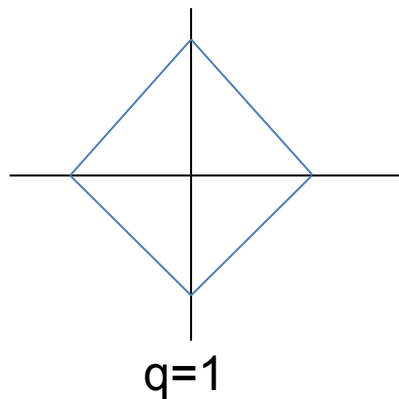
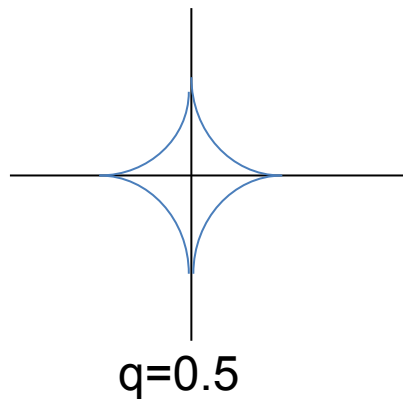
一般的な正則化項

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(y_i - \langle \phi(\mathbf{x}_i), \mathbf{w} \rangle \right)^2 + \underbrace{\frac{\lambda}{2} \sum_{j=1}^K |w_j|^q}_{\text{正則化項}}$$

- $q=2$ のときがL2正則化
- $q=1$ のときはLASSO: 1ノルムによる正則化なので **L1正則化**と呼ぶ
- Least Absolute Shrinkage and Selection Operator
 - λ が十分大きいと、 w_j のいくつかは0になりやすい → スパースなモデル
- $q=0$ のときはL0正則化。解きにくい問題(上記2つと違い凸ではない)

- 制約 $\sum_{j=1}^K |w_j|^q \leq \eta$

のもとで、 $L(w)$ を最小化する、と考える。



L2正則化

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(y_i - \langle \phi(\mathbf{x}_i), \mathbf{w} \rangle \right)^2 + \underbrace{\frac{\lambda}{2} \mathbf{w}^T \mathbf{w}}_{\text{正則化項}}$$

最小化すると

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w}) = (\lambda \mathbf{I} + \phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^T \mathbf{y}$$

正則化項

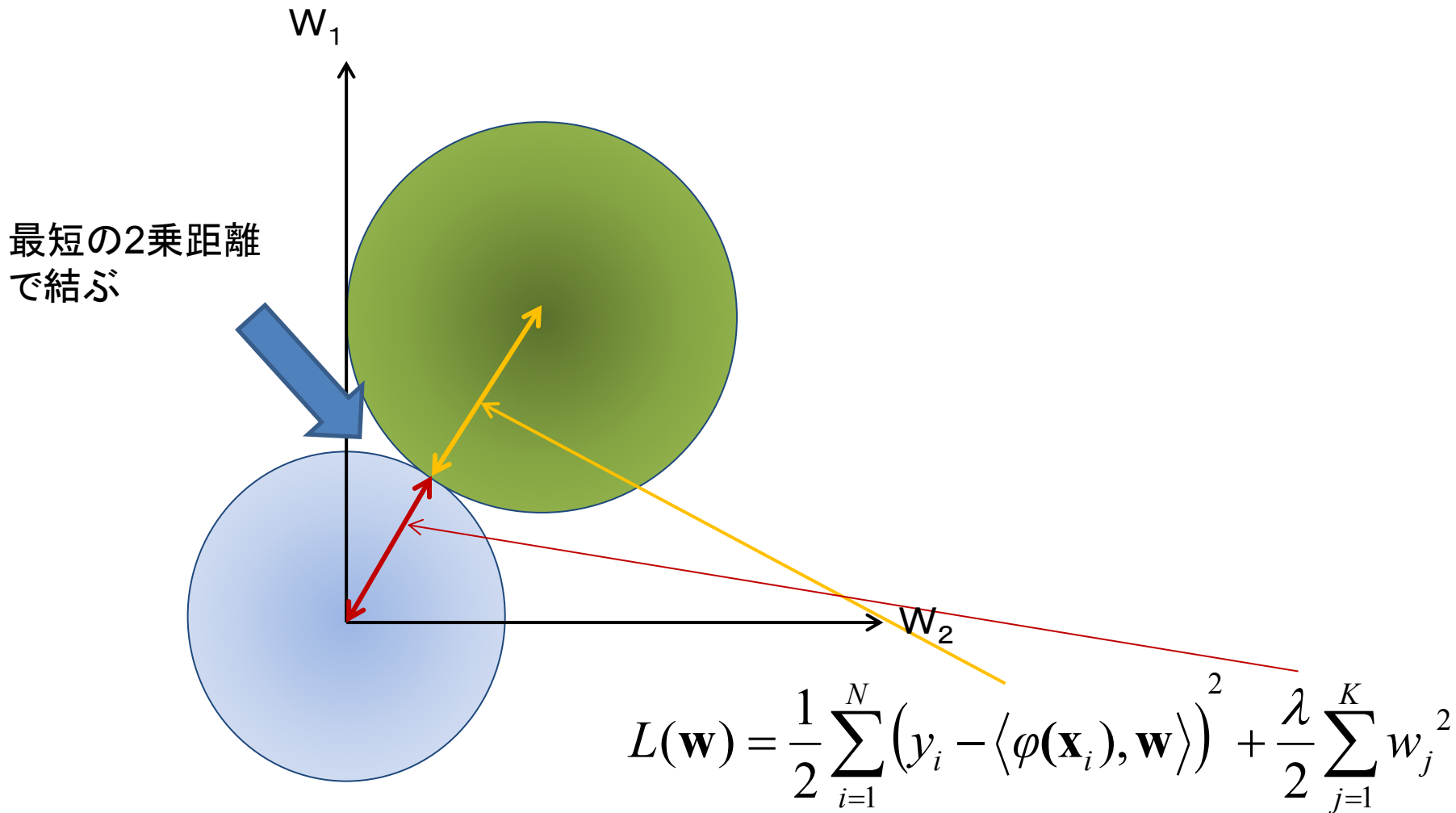
(\mathbf{w} の影響を小さくする効果)

\mathbf{W} の2ノルムによる正則化であるので、L2正則化と呼ぶ

➤ 最適な \mathbf{w} は $L(\mathbf{w})$ を微分して0とすれば上記のように解析的に閉じた式で求まる。

➤ これは $\phi(\mathbf{X})$ と λ の案配によって決まり、どの成分も強制的にゼロにしようという力は働かない

L2正則化のイメージ

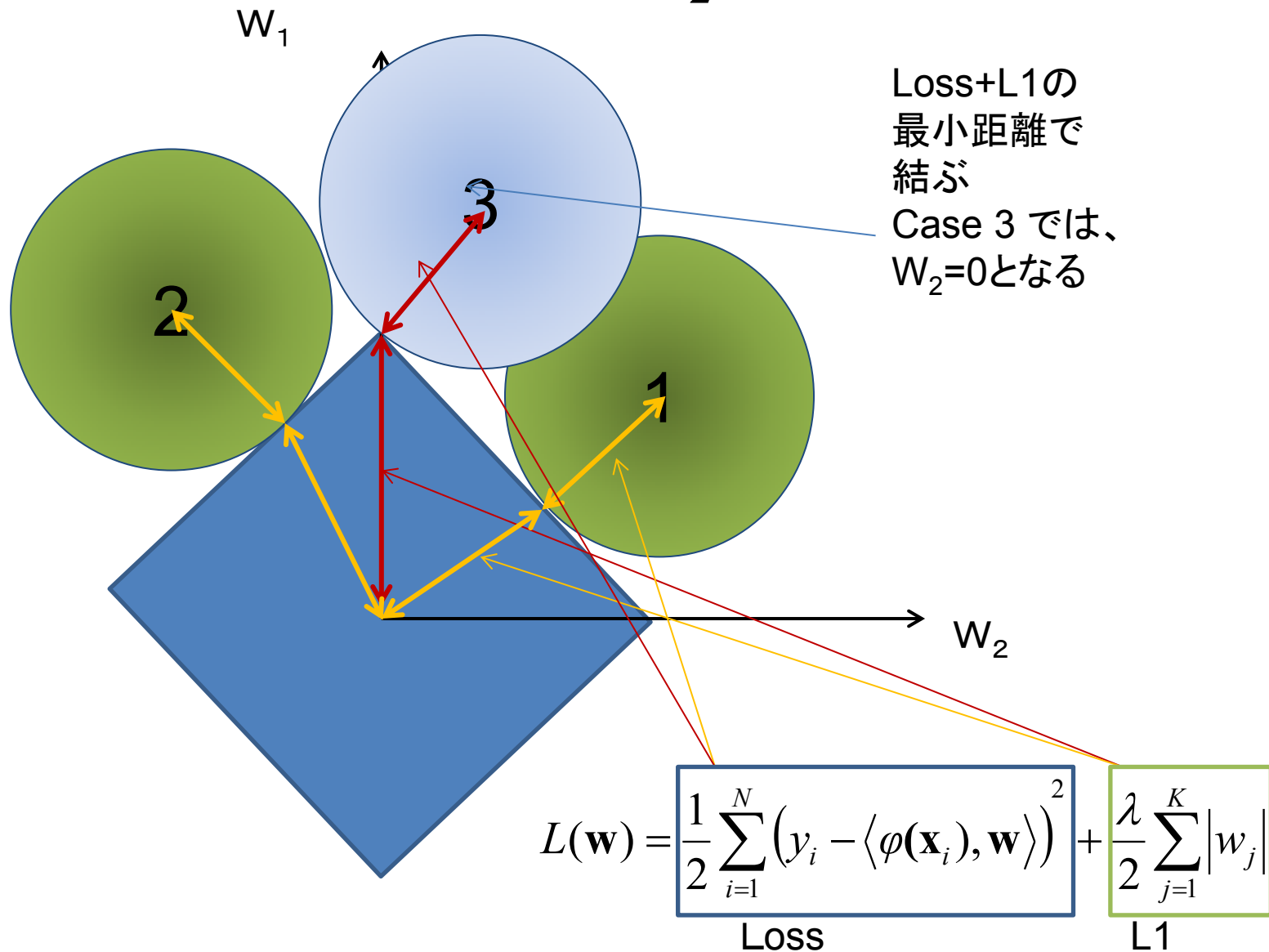


L1正則化

- L2正則化では w の最適値 \hat{w} を損失 L の微分で閉じた式で求められたが、L1正則化では $|w|$ が $w=0$ で微分できないので、ややこしくなる。
- L1正則化を行う逐次的な方法と
L1正則化が w の要素の多くをゼロ化する傾向を以下で説明する

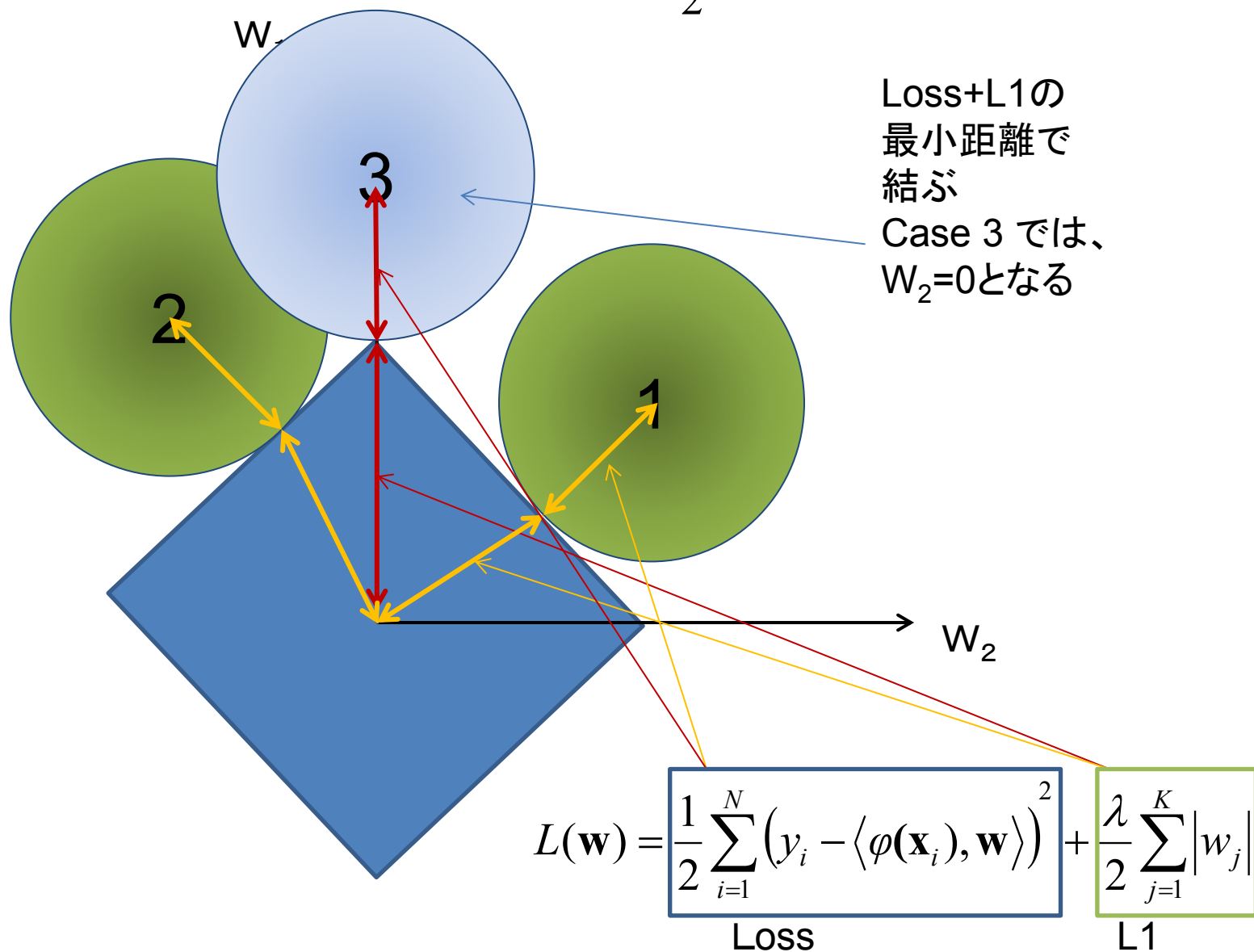
L1正則化イメージ: (1)

w_2 軸でのLossの微分=0として \tilde{w}_2 を求める



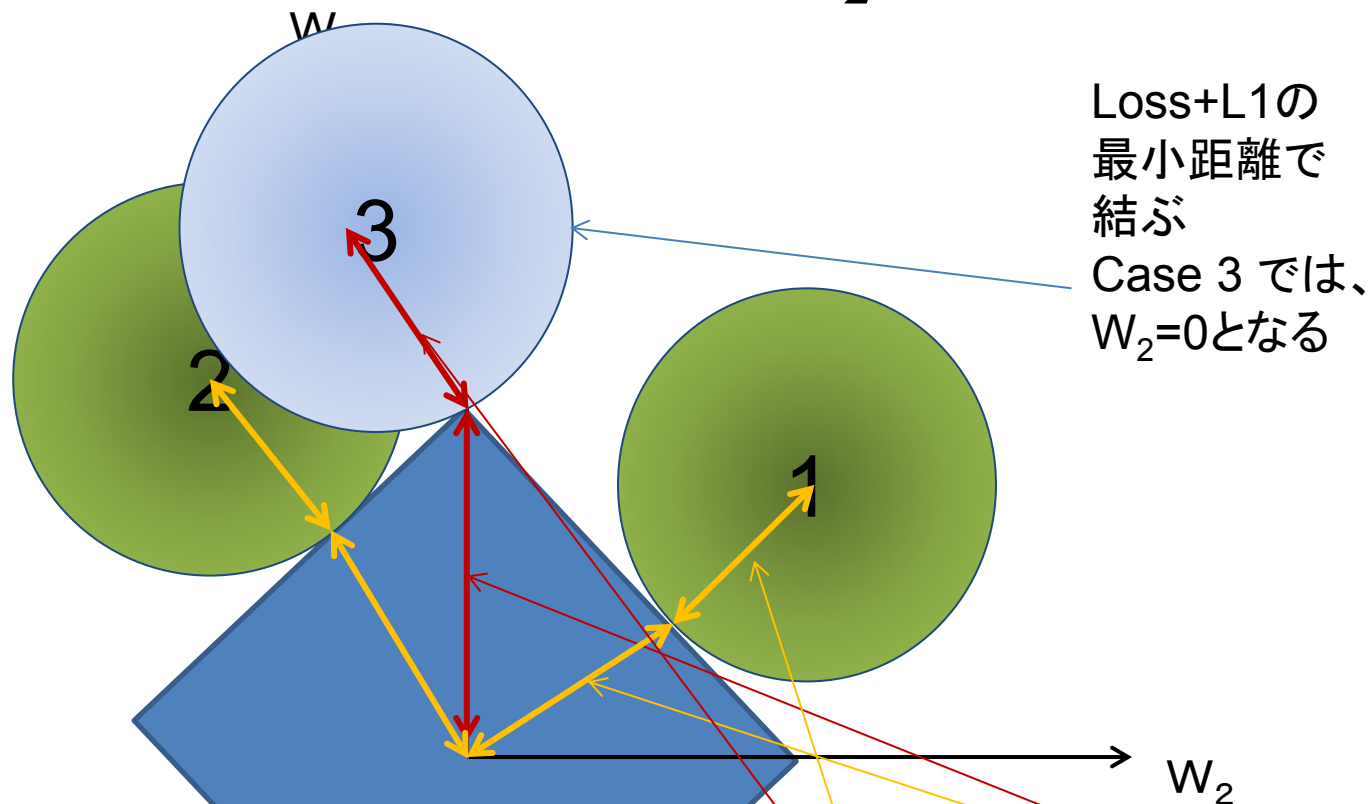
L1正則化イメージ: (2)

w_2 軸でのLossの微分=0として \tilde{w}_2 を求める



L1正則化イメージ: (3)

w_2 軸でのLossの微分=0として \tilde{w}_2 を求める



(1)(2)(3)で2本の赤い矢印線の長さの和が変わらない点に注目
以下でL1正則化に関してもう少し細かく議論する。

$$L(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{i=1}^N (y_i - \langle \phi(\mathbf{x}_i), \mathbf{w} \rangle)^2}_{\text{Loss}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^K |w_j|}_{\text{L1}}$$

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(y_i - \langle \phi(\mathbf{x}_i), \mathbf{w} \rangle \right)^2 + \frac{\lambda}{2} \sum_{j=1}^K |w_j| \quad (L1-10)$$

- ある次元 d に着目して $L(\mathbf{w})$ を最小化するような w_d を求める。
- これを各次元について繰り返し、 $L(\mathbf{w})$ の最小化を図る。 w_d について $L(\mathbf{w})$ を書き直すと

$$\begin{aligned} L(\mathbf{w}) &= \frac{1}{2} \sum_{i=1}^N \left(y_i - \phi_d(\mathbf{x}_i)w_d - \sum_{j \neq d} \phi_j(\mathbf{x}_i)w_j \right)^2 + \frac{\lambda}{2} \left(|w_d| + \sum_{j \neq d} |w_j| \right) \\ &= \text{Loss}(\mathbf{w}) \quad + \quad L1(\mathbf{w}) \quad (L1-20) \end{aligned}$$

- $\frac{\partial L(\mathbf{w})}{\partial w_d} = 0$ とおき w_d の最適値を求めたいが絶対値を含む第2項 $L1(w)$ が微分できないので、ひとまず $\text{Loss}(\mathbf{w})$ を微分して0とおくと

$$\begin{aligned}
\frac{\partial \text{Loss}(\mathbf{w})}{\partial w_d} &= \frac{\partial}{\partial w_d} \left(\frac{1}{2} \sum_{i=1}^N \left(y_i - \varphi_d(\mathbf{x}_i) w_d - \sum_{j \neq d} \varphi_j(\mathbf{x}_i) w_j \right)^2 \right) \\
&= \sum_{i=1}^N -\varphi_d(\mathbf{x}_i) \left(y_i - \varphi_d(\mathbf{x}_i) w_d - \sum_{j \neq d} \varphi_j(\mathbf{x}_i) w_j \right) = 0 \text{ の解を } \tilde{w}_d \text{ とする} \\
\Rightarrow \quad \tilde{w}_d &= \frac{\sum_{i=1}^N \varphi_d(\mathbf{x}_i) \left(y_i - \sum_{j \neq d} \varphi_j(\mathbf{x}_i) w_j \right)}{\sum_{i=1}^N \varphi_d(\mathbf{x}_i)^2}
\end{aligned}$$

これを用いて $L(\mathbf{w})$ を書き換える。ただし、 w_d に関係しないところは当面定数と見なせるので、無視した。

$$\begin{aligned}
2L(\tilde{w}_d) &= \sum_{i=1}^N \left(\varphi_d(\mathbf{x}_i)^2 w_d^2 - 2\varphi_d(\mathbf{x}_i) w_d \left(y_i - \sum_{j \neq d} \varphi_j(\mathbf{x}_i) w_j \right) \right) + \lambda |w_d| + \text{Const} \\
&= w_d^2 \sum_{i=1}^N \varphi_d(\mathbf{x}_i)^2 - 2w_d \sum_{i=1}^N \left(\varphi_d(\mathbf{x}_i)^2 \right) \sum_{i=1}^N \left(\varphi_d(\mathbf{x}_i) w_d \left(y_i - \sum_{j \neq d} \varphi_j(\mathbf{x}_i) w_j \right) \right) \bigg/ \sum_{i=1}^N \left(\varphi_d(\mathbf{x}_i)^2 \right) + \lambda |w_d| + \text{Const} \\
&= w_d^2 \sum_{i=1}^N \varphi_d(\mathbf{x}_i)^2 - 2w_d \tilde{w}_d \sum_{i=1}^N \varphi_d(\mathbf{x}_i)^2 + \lambda |w_d| + \text{Const}
\end{aligned}$$

ここで $\gamma = \frac{\lambda}{2 \sum_{i=1}^N \varphi_d(\mathbf{x}_i)^2}$ とおくと、 $L(w_d) = \frac{1}{2}(w_d^2 - w_d \tilde{w}_d) + \gamma |w_d| + \text{const}$

$$\frac{\partial L(w_d)}{\partial w_d} = \begin{cases} w_d - \tilde{w}_d + \gamma & w_d > 0 \\ w_d - \tilde{w}_d - \gamma & w_d < 0 \\ \text{undefined} & w_d = 0 \end{cases} \Rightarrow \frac{\partial L(w_d)}{\partial w_d} = 0 \text{ なる } w_d \text{ を探す}$$

case 1 $\tilde{w}_d - \gamma > 0$ なら $w_d > 0$ なので $w_d = \tilde{w}_d - \gamma$

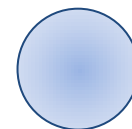
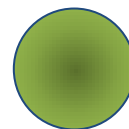
case 2 $\tilde{w}_d + \gamma < 0$ なら $w_d < 0$ なので $w_d = \tilde{w}_d + \gamma$

case 3 $-\gamma < \tilde{w}_d < \gamma$ なら $w_d = 0$ なぜなら

$w_d > 0$ だと $\tilde{w}_d - \gamma > 0$ すなわち $\tilde{w}_d > \gamma$ 矛盾

$w_d < 0$ だと $\tilde{w}_d + \gamma < 0$ すなわち $\tilde{w}_d < -\gamma$ 矛盾

case 3により $-\gamma < \tilde{w}_d < \gamma$ すなわち $\text{Loss}(\mathbf{w})$ の w_d の解 \tilde{w}_d がゼロに近づくと $w_d = 0$ になりゼロ化(スパース化) される力が働く



W全体の正則化

[step 1] w の各要素を適当な値に初期化

[step 2] w の各要素の値 $w_k (k=1, \dots, K)$ が収束するまで以下 step 3, 4, 5 を繰り返す

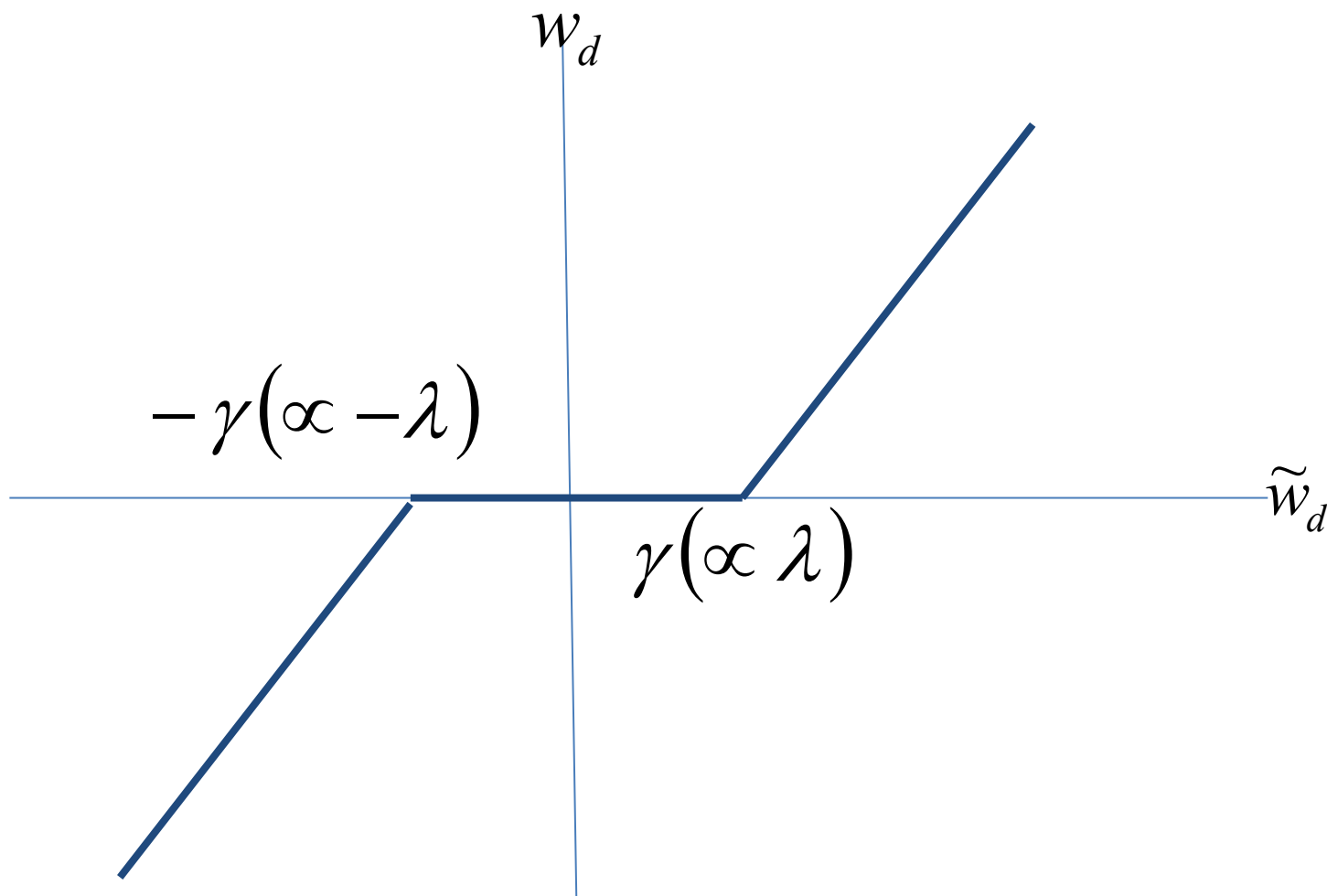
[step 3] $k=1, \dots, K$ で step 4, step 5 を繰り返す

[step 4] $w_j (j \neq k)$ を用いて case 1, 2, 3 にしたがって w_j を計算してゼロ化

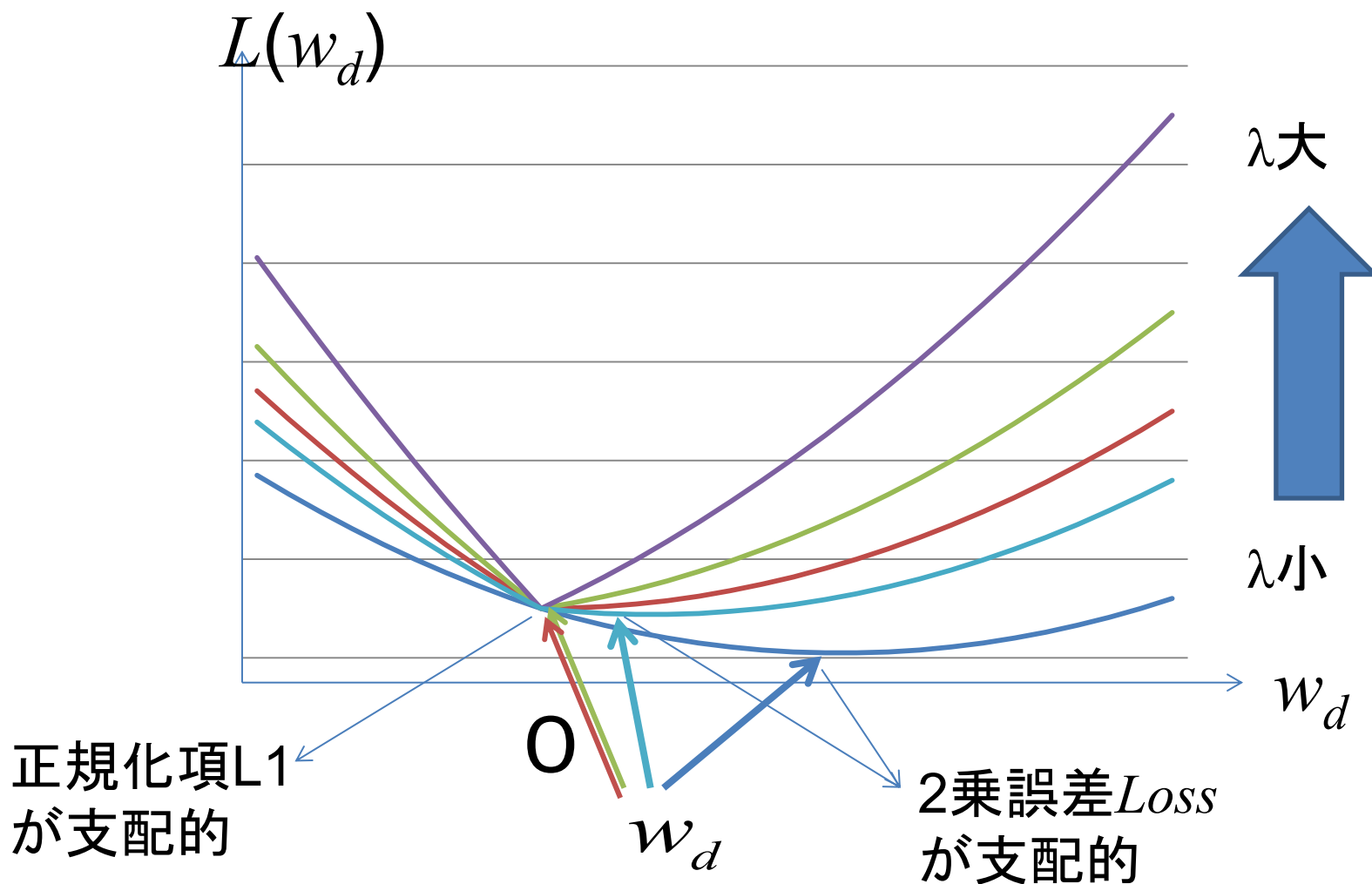
[step 5] w_k を更新

[step 6] 収束したら w の最終結果とする

w_d のゼロ化のイメージ



L1正則化が支配的になり \hat{w}_d をゼロ化する様子を下図で例示する



正則化項のBayes的解釈

➤ Bayesでは事後確率は

観測データの確率 × 事前確率

➤ 事後確率を最大化するパラメタ η を求めたい

$$\hat{\eta} = \arg \max_{\eta} P(X | \eta) P(\eta | \alpha) \quad \alpha \text{は事前分布のハイパーパラメタ}$$

➤ ここで対数尤度にとしてみると、次のように解釈できる

$$\hat{\eta} = \arg \max_{\eta} (\log P(X | \eta) + \log P(\eta | \alpha))$$

損失関数

正則化項

例: 事前分布、事後分布とも正規分布

$$\mathbf{y} = (y_1, \dots, y_N)^T \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_K \end{bmatrix}$$

$$y = \varphi(\mathbf{x})\mathbf{w} + \varepsilon \quad \varepsilon = N(0,1)$$

$$\log \prod_i p(y_i | \mathbf{x}_i, \mathbf{w}, 1) = \sum_i \log N(y_i | \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle, 1) \propto \sum_i - (y_i - \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle)^2 / 2$$

事前分布 $p(\mathbf{w} | \alpha, \gamma)$ も同様にすると

$$\log p(\mathbf{w} | \alpha, \gamma) \propto -(\mathbf{w} - \alpha)^T (\mathbf{w} - \alpha) / 2$$

$$\Rightarrow \arg \min_{\mathbf{w}} \left(\log \prod_i p(y_i | \mathbf{x}_i, \mathbf{w}, \beta) + \log p(\mathbf{w} | \alpha, \gamma) \right)$$

$$\arg \min_{\mathbf{w}} \left(\frac{1}{2} \sum_i (y_i - \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle)^2 + \frac{1}{2} (\mathbf{w} - \alpha)^T (\mathbf{w} - \alpha) \right)$$

ここで、 $\alpha = 0$, 事前分布の重みを λ とすると

$$\Rightarrow \arg \max_{\mathbf{w}} \left(\frac{1}{2} \sum_i (y_i - \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right) \quad \text{L2ノルムによる正則化項}$$

事前分布の \mathbf{w} の
分散: λ^{-1} とも見
える。

例：事前分布がLaplace分布、事後分布が正規分布

$$y = \varphi(\mathbf{x})\mathbf{w} + \varepsilon \quad \varepsilon = N(0,1)$$

$$\log \prod_i p(y_i | \mathbf{x}_i, \mathbf{w}, 1) = \sum_i \log N(y_i | \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle, 1) \propto \sum_i - (y_i - \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle)^2 / 2$$

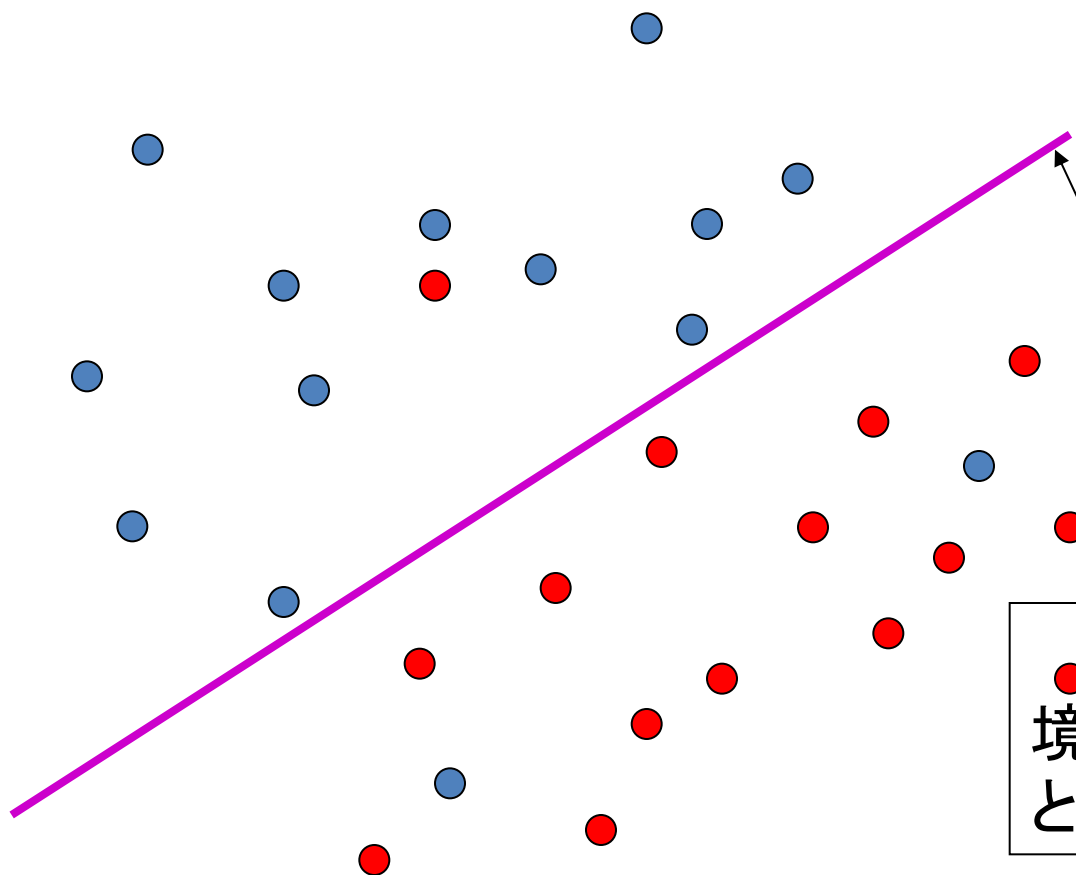
事前分布は期待値0のLaplace分布 $p(\mathbf{w} | \lambda) = \frac{\lambda}{4} \exp\left(-\frac{\lambda|\mathbf{w}|}{2}\right)$ も同様にすると

$$\log p(\mathbf{w} | \lambda) \propto -\frac{\lambda|\mathbf{w}|}{2}$$

$$\begin{aligned} \Rightarrow \arg \min_{\mathbf{w}} & \left(\log \prod_i p(y_i | \mathbf{x}_i, \mathbf{w}, \beta) + \log p(\mathbf{w} | \lambda) \right) \\ & = \arg \min_{\mathbf{w}} \left(\frac{1}{2} \sum_i (y_i - \langle \varphi(\mathbf{x}_i), \mathbf{w} \rangle)^2 + \frac{\lambda}{2} |\mathbf{w}| \right) \quad \text{L1ノルムによる正則化項} \end{aligned}$$

- 以上、述べてきた線形回帰によるモデル化は、生成モデル
- 当然、線形の識別モデルもある。次以降は線形識別モデルの話

線形識別



● と ● の領域の
境界面を線形関数
として求める

線形識別

- データ: $\mathbf{x} = [x_1, x_2, \dots, x_M]^T$
- \mathbf{x} がいくつかのクラス(あるいはカテゴリー): C_k のどれかに属する。
 - 例: 新聞記事が「政治」「経済」「スポーツ」「芸能」「社会」などのクラスのどれかに属する場合。この場合、データ: \mathbf{x} は例えば、記事に現れる単語の集合、など。
- データ: \mathbf{x} がK個のクラスの各々に属するかどうかの判定は(−1 = 属さない, 1 = 属する)の2値を要素とするK次元ベクトル: $\mathbf{y}_i = (-1, 1, -1, \dots, 1)$ で表される。
 - ただし、1つのクラスに属するか属さないかだけを識別すの場合は2クラス分類という。当然、 $y_i = -1$ or $y_i = 1$
- この属するか否かの判断をする式が線形の場合を線形識別という。

➤ 線形識別の関数

$$y(\mathbf{x}) = \langle \mathbf{x}, \mathbf{w} \rangle + w_0$$

あるいは $\tilde{\mathbf{x}} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}, \tilde{\mathbf{w}} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$ とおくなら $y(\mathbf{x}) = \langle \tilde{\mathbf{x}}, \tilde{\mathbf{w}} \rangle$

➤ 一般化線形識別の関数は以下


$$y(\mathbf{x}) = f(\langle \mathbf{x}, \mathbf{w} \rangle + w_0) \quad f \text{ は非線形でもよい}$$

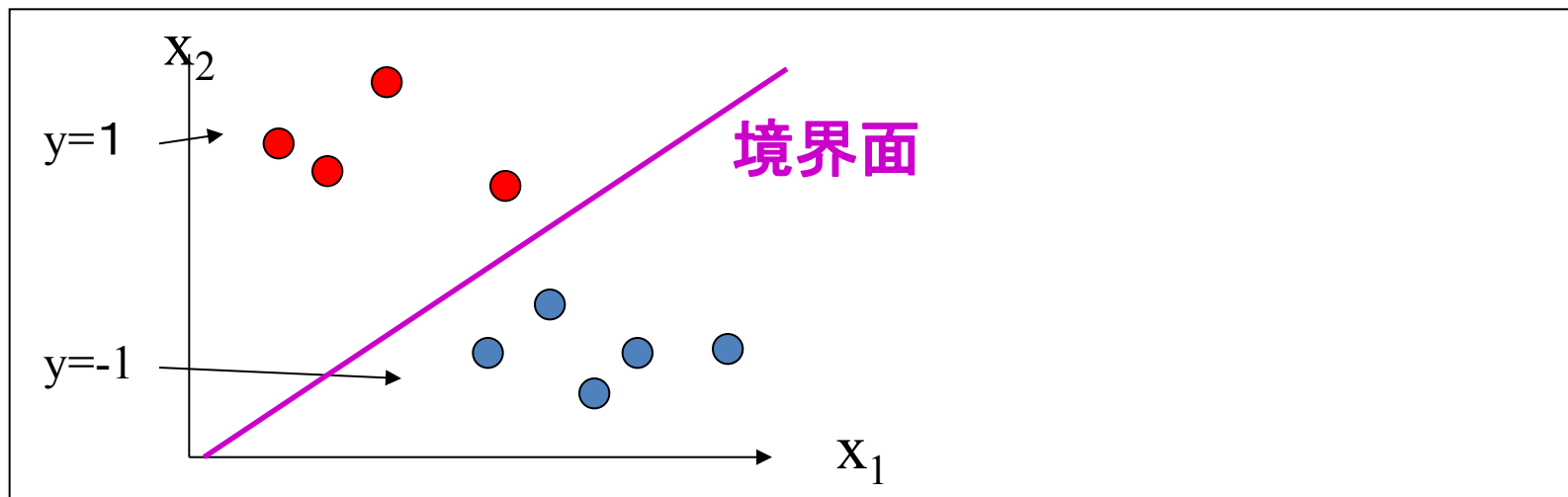
➤ 2クラス分類

➤ クラス C_1 に属するか $C_2 (= \text{not} C_1)$ に属するかは、次の通り

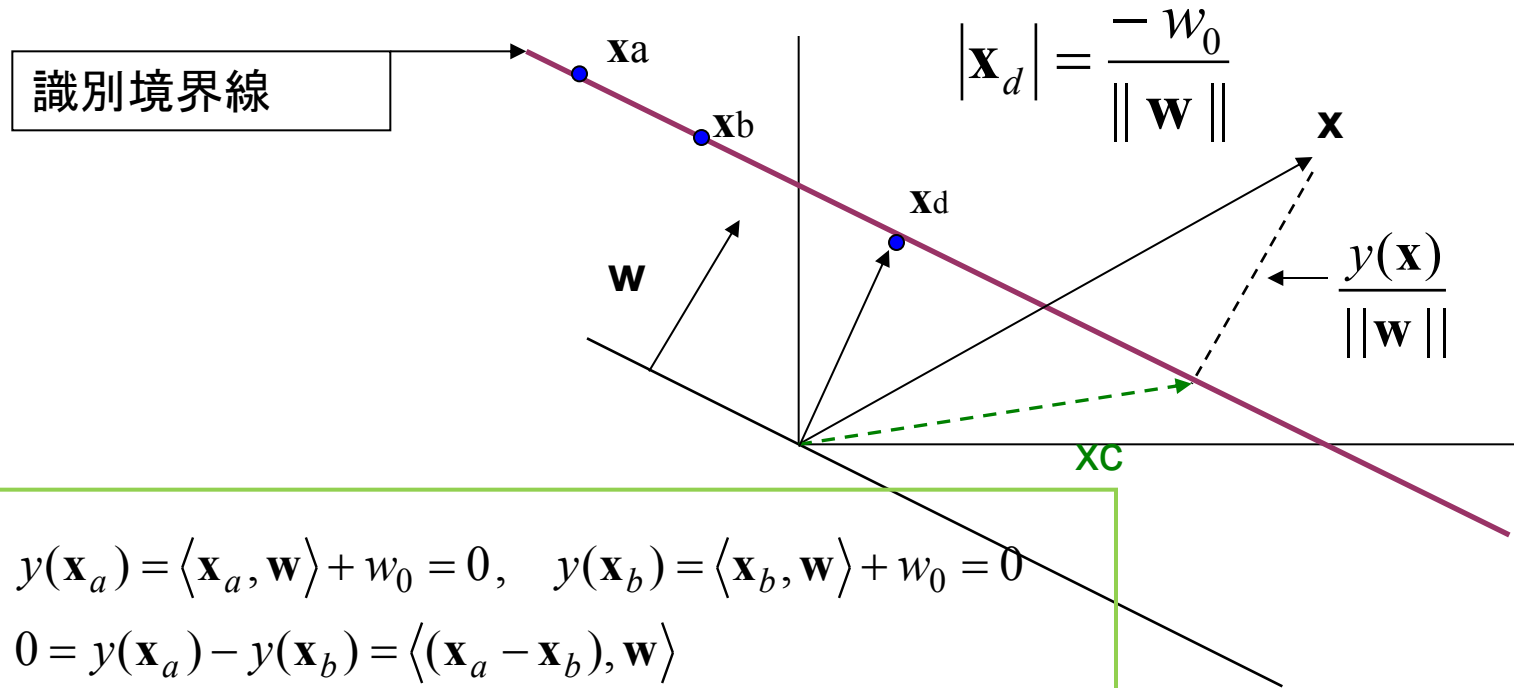
➤ if $y(\mathbf{x}) \geq 0$ then データ: \mathbf{x} は C_1 に属する
otherwise データ: \mathbf{x} は C_2 に属する
(すなわち C_1 に属さない)

2値分類の直観的説明

- $y=\{-1,1\}$ 、 x は2次元とする。(下図を参照)
- $\{y,x\}$ を教師データとして、2乗誤差の最小化を行って正規方程式を求めると、下図の  のようなクラスを分類する分離平面が得られる。



線形識別関数の幾何学的解釈



原点から識別境界線への垂線の交点を \mathbf{x}_d とおく。

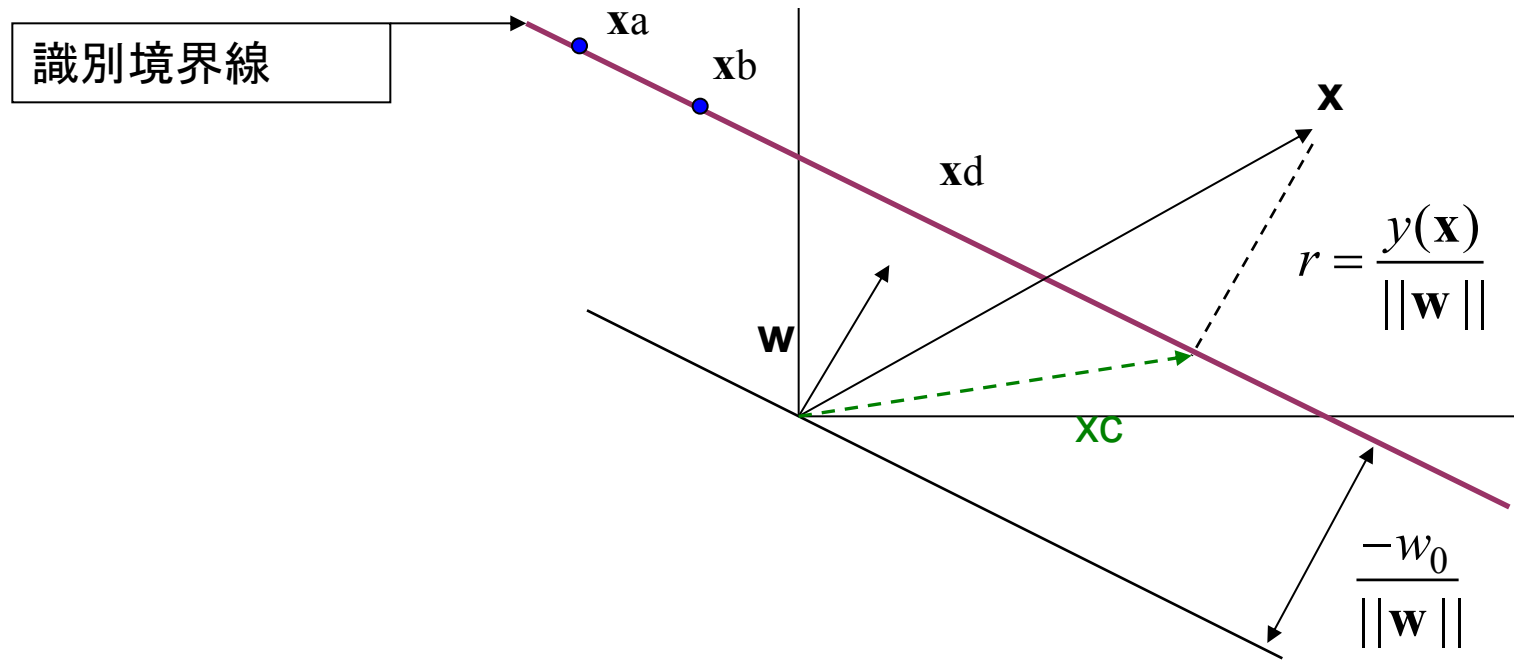
$$0 = y(\mathbf{x}_d) = \langle \mathbf{x}_d, \mathbf{w} \rangle + w_0$$

\mathbf{x}_d は \mathbf{w} に並行で横ベクトルだから、 $\langle \mathbf{x}_d, \mathbf{w} \rangle = \|\mathbf{x}_d\| \cdot \|\mathbf{w}\|$

これを上式に代入して整理すると

$$\langle \mathbf{x}_d, \mathbf{w} \rangle + w_0 = \|\mathbf{x}_d\| \cdot \|\mathbf{w}\| + w_0 = 0 \quad \Rightarrow \quad |\mathbf{x}_d| = -\frac{w_0}{\|\mathbf{w}\|}$$

線形識別関数の幾何学的解釈



$$\mathbf{x} = \mathbf{x}_c + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \quad \text{両辺と } \mathbf{w} \text{ の内積をとり、} w_0 \text{ を足すと}$$

$$y(\mathbf{x}) = \langle \mathbf{x}, \mathbf{w} \rangle + w_0 = \langle \mathbf{x}_c, \mathbf{w} \rangle + w_0 + r \frac{\langle \mathbf{w}, \mathbf{w} \rangle}{\|\mathbf{w}\|} = y(\mathbf{x}_c) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|}$$

$$y(\mathbf{x}_c) = 0 \quad \text{だから} \quad r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$$

wの計算方法:2クラス分類の場合

- . クラス C_1, C_2 の境界が $y(\mathbf{x}) = \langle \tilde{\mathbf{x}}, \tilde{\mathbf{w}} \rangle$ で書けるとする
- すると新規のデータ: \mathbf{x} は $y(\tilde{\mathbf{x}})$ が正ならクラス C_1 に, 負なら C_2 属する

N 個の教師データ $\{\tilde{\mathbf{x}}_n, y_n\} (n=1, N)$ があったとき

- . ただしクラス1なら $y_n = 1$, 0なら $y_n = -1$

$$\tilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_N^T \end{pmatrix} \quad \mathbf{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \quad \tilde{\mathbf{X}}\tilde{\mathbf{W}} = \begin{pmatrix} \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}} \rangle \\ \vdots \\ \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}} \rangle \end{pmatrix}$$

- すると、観測データ(教師データ)において個々のクラスに分類されたか否かの観点からの2乗誤差は次式となる

$$E(\tilde{\mathbf{W}}) = (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})$$

- もう少し詳しく書くと

$$\begin{aligned} & (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}) = \\ & = \left(\langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}} \rangle - y_1 \quad \cdots \quad \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}} \rangle - y_N \right) \begin{pmatrix} \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}} \rangle - y_1 \\ \vdots \\ \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}} \rangle - y_N \end{pmatrix} \\ & = \left(\langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}} \rangle - y_1 \right)^2 + \cdots + \left(\langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}} \rangle - y_N \right)^2 \end{aligned}$$

$$E(\tilde{\mathbf{W}}) = (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})$$

➤ これを最小化する $\tilde{\mathbf{W}}$ は $\tilde{\mathbf{W}}$ で微分して0とおけば、線形回帰のときと同様の計算により求まる。

➤ 微分は次式：

$$\frac{\partial \mathbf{A}^T \mathbf{A}}{\partial \mathbf{W}} = 2 \frac{\partial \mathbf{A}^T}{\partial \mathbf{W}} \mathbf{A} \rightarrow \mathbf{A} = (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}) \rightarrow = 2\tilde{\mathbf{X}}^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})$$

$$\begin{aligned} \frac{\partial E(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{W}}} &= \tilde{\mathbf{X}}^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}) = 0 \\ \Rightarrow \quad \tilde{\mathbf{W}} &= (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y} \end{aligned}$$

- 新規のデータ \mathbf{x}_{new} に対する予測を行う $y(\mathbf{x}_{\text{new}})$ も求まる。

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

$$\mathbf{y}(\tilde{\mathbf{x}}_{\text{new}}) = \begin{bmatrix} y_1(\tilde{\mathbf{x}}_{\text{new}}) \\ \vdots \\ y_K(\tilde{\mathbf{x}}_{\text{new}}) \end{bmatrix} = \tilde{\mathbf{x}}_{\text{new}} \tilde{\mathbf{W}} = \tilde{\mathbf{x}}_{\text{new}} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

- $\mathbf{y}(\mathbf{x}_{\text{new}})$ が大きいほどクラス C_1 に属する可能性が高い。

wの計算方法：多クラス分類の場合

- クラス C_k が線形識別モデル $y_k(\mathbf{x}) = \tilde{\mathbf{x}}\tilde{\mathbf{w}}_k$ で書けるとする。
- すると新規のデータ: \mathbf{x} は $y_k(\tilde{\mathbf{x}})$ が最大の k のクラス C_k に属する

$$y_k(\mathbf{x}) \text{を} K \text{個並べたベクトル} \mathbf{y} = [y_1(\mathbf{x}) \quad \cdots \quad y_K(\mathbf{x})]^T \\ = \left(\langle \tilde{\mathbf{x}}, \tilde{\mathbf{w}}_1 \rangle \quad \cdots \quad \langle \tilde{\mathbf{x}}, \tilde{\mathbf{w}}_K \rangle \right) = \tilde{\mathbf{x}}\tilde{\mathbf{W}}$$

- N 個の教師データ $\{\tilde{\mathbf{x}}_n, \mathbf{y}_n\} (n=1, \dots, N)$ があったとき
注 $\tilde{\mathbf{x}}_n$ は K 個のクラス内の複数個に属することもあるなら
 \mathbf{y}_n は K 次元ベクトル $(-1, 1, -1, \dots, 1)$ のような形。

$$\tilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_N^T \end{pmatrix} \quad \mathbf{Y} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_N \end{pmatrix} \quad \tilde{\mathbf{X}}\tilde{\mathbf{W}} = \begin{pmatrix} \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_1 \rangle & \cdots & \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_K \rangle \\ \vdots & \ddots & \vdots \\ \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_1 \rangle & \cdots & \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_K \rangle \end{pmatrix}$$

- すると、観測データ(教師データ)において個々のクラスに分類されたか否かの観点からの2乗誤差は次式となる

$$E(\tilde{\mathbf{W}}) = Tr\left\{(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})\right\}$$

- もう少し詳しく書くと

$$\begin{aligned} & (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}) = \\ & = \begin{pmatrix} \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_1 \rangle - y_{11} & & \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_1 \rangle - y_{N1} \\ \vdots & \ddots & \\ \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_K \rangle - y_{1K} & & \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_K \rangle - y_{NK} \end{pmatrix} \begin{pmatrix} \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_1 \rangle - y_{11} & & \langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_K \rangle - y_{1K} \\ & \ddots & \\ \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_1 \rangle - y_{N1} & & \langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_K \rangle - y_{NK} \end{pmatrix} \\ & \therefore Tr\left((\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})\right) \\ & = (\langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_1 \rangle - y_{11})^2 + \cdots + (\langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_1 \rangle - y_{N1})^2 \\ & \quad + \cdots + (\langle \tilde{\mathbf{x}}_1, \tilde{\mathbf{w}}_K \rangle - y_{1K})^2 + \cdots + (\langle \tilde{\mathbf{x}}_N, \tilde{\mathbf{w}}_K \rangle - y_{NK})^2 \end{aligned}$$

$$E(\tilde{\mathbf{W}}) = \text{Tr}\left\{\left(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}\right)^T \left(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}\right)\right\}$$

➤ これを最小化する $\tilde{\mathbf{W}}$ は $\tilde{\mathbf{W}}$ で微分して0とおけば、線形回帰のときと同様の計算により求まる。

➤ Tr の微分は次式：

$$\frac{\partial \text{Tr}(\mathbf{A}^T \mathbf{A})}{\partial \mathbf{W}} = 2 \frac{\partial \mathbf{A}^T}{\partial \mathbf{W}} \mathbf{A} \rightarrow \mathbf{A} = (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}) \rightarrow = 2\tilde{\mathbf{X}}^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y})$$

$$\frac{\partial E(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{W}}} = \tilde{\mathbf{X}}^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{Y}) = 0$$

$$\Rightarrow \tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

- 新規のデータ \mathbf{x}_{new} に対する予測を行う $y(\mathbf{x}_{\text{new}})$ も求まる。

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

$$\mathbf{y}(\tilde{\mathbf{x}}_{\text{new}}) = \begin{bmatrix} y_1(\tilde{\mathbf{x}}_{\text{new}}) \\ \vdots \\ y_K(\tilde{\mathbf{x}}_{\text{new}}) \end{bmatrix} = \tilde{\mathbf{x}}_{\text{new}} \tilde{\mathbf{W}} = \tilde{\mathbf{x}}_{\text{new}} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{Y}$$

- $y_i(\mathbf{x}_{\text{new}})$ が大きいほどそのクラス i に属する可能性が高い。

もちろん、 $y_i(\mathbf{x}_{\text{new}})$ が最大となる i のクラスに属すると考えるのが自然。だが。。。

生成モデルを利用した識別

- 識別はベイズ統計的には次式

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{p(\mathbf{x})}$$

- N個のデータ: \mathbf{x}_k ($k=1, \dots, N$)があるクラスに属するかどうかの判定は(0=属さない, 1=属する)の2値を要素とするN個のK次元ベクトル: $\mathbf{y} = (0, 1, 0, \dots, 1)$ で表される。
 - 以下のベイズ統計による分類では、属さない場合を-1ではなく0とすることに注意。
- 以下ではベイズ統計による2クラス分類をする場合に事後確率について考える。

Logistic sigmoid function

➤ クラス C_1 の事後分布は次式(s-1)

$$\begin{aligned} p(C_1 | \mathbf{x}) &= \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \quad \text{-(s-1) logistic sigmoid function} \end{aligned}$$

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

$$\sigma(-a) = 1 - \sigma(a) \quad a = \log \frac{\sigma}{1 - \sigma}$$

$$\frac{d\sigma}{da} = \frac{\exp(-a)}{(1 + \exp(-a))^2} = \frac{1}{1 + \exp(-a)} \cdot \frac{\exp(-a)}{1 + \exp(-a)} = \sigma(1 - \sigma)$$

クラス C_1, C_2 が共分散 Σ が等しい2つの正規分布の場合の事後確率 $p(C_i|\mathbf{x})$

➤ 式(s-1)によって以下のように導ける。

$$p(\mathbf{x} | C_i) = \frac{1}{(2\pi)^{K/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_i)^T \Sigma^{-1}(\mathbf{x} - \mu_i)\right\}$$

$$\begin{aligned} & \log \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} \\ &= \frac{\log\left((2\pi)^{K/2} |\Sigma|^{1/2}\right)}{\log\left((2\pi)^{K/2} |\Sigma|^{1/2}\right)} \cdot \left(-\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) - \frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma^{-1}(\mathbf{x} - \mu_2)\right) + \log \frac{p(C_1)}{p(C_2)} \\ &= \left(\frac{1}{2}(\mathbf{x}^T \Sigma^{-1} \mu_1 + \mu_1^T \Sigma^{-1} \mathbf{x} - \mu_1^T \Sigma^{-1} \mu_1) - \frac{1}{2}(\mathbf{x}^T \Sigma^{-1} \mu_2 + \mu_2^T \Sigma^{-1} \mathbf{x} - \mu_2^T \Sigma^{-1} \mu_2)\right) + \log \frac{p(C_1)}{p(C_2)} \\ &= \left(\mathbf{x}^T \Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2\right) + \log \frac{p(C_1)}{p(C_2)} \end{aligned}$$

Σ が2つのクラスで等しいことにとってキャンセルしていることに注意。等しくないともう少し複雑。

クラス C_1, C_2 が共分散 Σ が等しい2つの正規分布の場合の事後確率 $p(C_1|\mathbf{x})$

$$p(\mathbf{x} | C_i) = \frac{1}{(2\pi)^{K/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu_i)^T \Sigma^{-1}(\mathbf{x} - \mu_i)\right\}$$

$$\log \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \left(\mathbf{x}^T \Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 \right) + \log \frac{p(C_1)}{p(C_2)}$$

$$\Rightarrow \text{事後確率 : } p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = \frac{1}{1 + \exp(-(\mathbf{w}^T \mathbf{x} + w_0))}$$

$$\text{where } \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2)$$

$$w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \log \frac{p(C_1)}{p(C_2)}$$

Σ が2つのクラスで等しいことによってキャンセルしていることに注意。等しくないともう少し複雑。

次に Maximum likelihood solution (つまり \mathbf{w}, w_0)を求める。これによって、各クラスの事後確率が求まる

➤ ここで各クラスの事前確率が以下だったとする

$$\underline{p(C_1) = \pi \quad p(C_2) = 1 - \pi}$$

このとき観測データ \mathbf{x}_n が C_1 に属するとき $t_n = 1$ とし

$$p(\mathbf{x}_n, C_1) = p(C_1)p(\mathbf{x}_n | C_1) = \pi N(\mathbf{x}_n | \mu_1, \Sigma)$$

観測データ \mathbf{x}_n が C_2 に属するとき $t_n = 0$ とし

$$p(\mathbf{x}_n, C_2) = p(C_2)p(\mathbf{x}_n | C_2) = (1 - \pi)N(\mathbf{x}_n | \mu_2, \Sigma)$$

ここで $likelihood$ は次式 観測データは N 個あることを思い出そう

$$p(\mathbf{t} | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi N(\mathbf{x}_n | \mu_1, \Sigma)]^{t_n} [(1 - \pi)N(\mathbf{x}_n | \mu_2, \Sigma)]^{1-t_n}$$

$$where \quad \mathbf{t} = (t_1, \dots, t_N)^T \quad \dots (s-10)$$

➤(s-10)のlogすなわち log likelihood function を最大化することが目標

➤まず、最大化する π を求める。

➤(s-10)のlogの π に関する部分は次式(s-20) $\log p(\pi)$

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

$$\frac{\partial \log p(\pi)}{\partial \pi} = 0 \quad \Rightarrow \quad \pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$

where N_1 はクラス C_1 に属するデータ数。
 N_2 はクラス C_2 に属するデータ数。

➤ 次に (s-10)の \log を最大化する μ_1 を求める。

➤ (s-10)の \log の μ_2 に関する部分は次式(s-30) $\log p(\mu_1)$

$$\log p(\mu_2) = \sum_{n=1}^N (1 - t_n) \log N(\mathbf{x}_n | \mu_2, \Sigma) = -\frac{1}{2} \sum_{n=1}^N (1 - t_n) (\mathbf{x}_n - \mu_2)^T \Sigma^{-1} (\mathbf{x}_n - \mu_2)$$

+ const

$$\frac{\partial \log p(\mu_2)}{\partial \mu_2} = 0 \quad \Rightarrow \quad \mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n$$

➤ 同様にして μ_1 も求めると

$$\log p(\mu_1) = \sum_{n=1}^N t_n \log N(\mathbf{x}_n | \mu_1, \Sigma) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \mu_1)^T \Sigma^{-1} (\mathbf{x}_n - \mu_1) + \text{const}$$

$$\frac{\partial \log p(\mu_1)}{\partial \mu_1} = 0 \quad \Rightarrow \quad \mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n$$

➤最後に (s-10)の \log を最大化する精度行列 $\Lambda = \Sigma^{-1}$ (C_1 と C_2 共分散) を求める。

➤(s-10)の \log の Σ に関する部分は次式(s-40) $\log p(\Sigma)$

$$\begin{aligned}\log p(\Lambda) &= \frac{1}{2} \sum_{n=1}^N t_n \log |\Lambda| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \mu_1)^T \Lambda (\mathbf{x}_n - \mu_1) \\ &\quad + \frac{1}{2} \sum_{n=1}^N (1-t_n) \log |\Lambda| - \frac{1}{2} \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \mu_2)^T \Lambda (\mathbf{x}_n - \mu_2) \\ &= \frac{N}{2} \log |\Lambda| - \frac{N}{2} \text{Tr}(\Lambda S) \quad \dots (s-40)\end{aligned}$$

➤ $\log p(\Lambda)$ を Λ で微分して0とおき、(s-10)の \log を最大化する $\Lambda = \Sigma^{-1}$ を求める。

➤まず第1項の微分は線形代数学の公式より

$$\frac{\partial \left(\frac{N}{2} \log |\Lambda| \right)}{\partial \Lambda} = \frac{N}{2} (\Lambda^{-1})^T = \frac{N}{2} (\Lambda^{-1}) \quad \dots (s-50)$$

$\because \Lambda$ が対称 $\Rightarrow \Lambda^{-1}$ が対称

$(s - 40)$ の S は次式

$$S = \frac{1}{N} \sum_{n \in C_1} (\mathbf{x}_n - \mu_1)(\mathbf{x}_n - \mu_1)^T + \frac{1}{N} \sum_{n \in C_2} (\mathbf{x}_n - \mu_2)(\mathbf{x}_n - \mu_2)^T$$

➤次は $Tr(\Lambda S)$ を Λ で微分して0とおき、 $\log p(\Lambda)$ を最大化する Λ を求める。

$$\frac{\partial Tr(\Lambda S)}{\partial \Lambda} = -S^T = -S \quad \dots (s - 60)$$

$$\therefore \frac{\partial \log p(\Lambda)}{\partial \Lambda} = \frac{N}{2} \Lambda^{-1} - \frac{N}{2} S = 0$$

$$\Rightarrow \Lambda^{-1} = \Sigma = S = \frac{1}{N} \sum_{n \in C_1} (\mathbf{x}_n - \mu_1)(\mathbf{x}_n - \mu_1)^T + \frac{1}{N} \sum_{n \in C_2} (\mathbf{x}_n - \mu_2)(\mathbf{x}_n - \mu_2)^T$$

➤このようにして、教師データ集合 $\{(\mathbf{x}_n, t_n) | n=1, \dots, N\}$ から $\mu_1, \mu_2, \Sigma^{-1}(=\Lambda), \pi$ が求まったので、これらを用いて定義される \mathbf{w}, w_0 も求まる。

➤未知データ \mathbf{x} がクラス C_1 に属する確率は

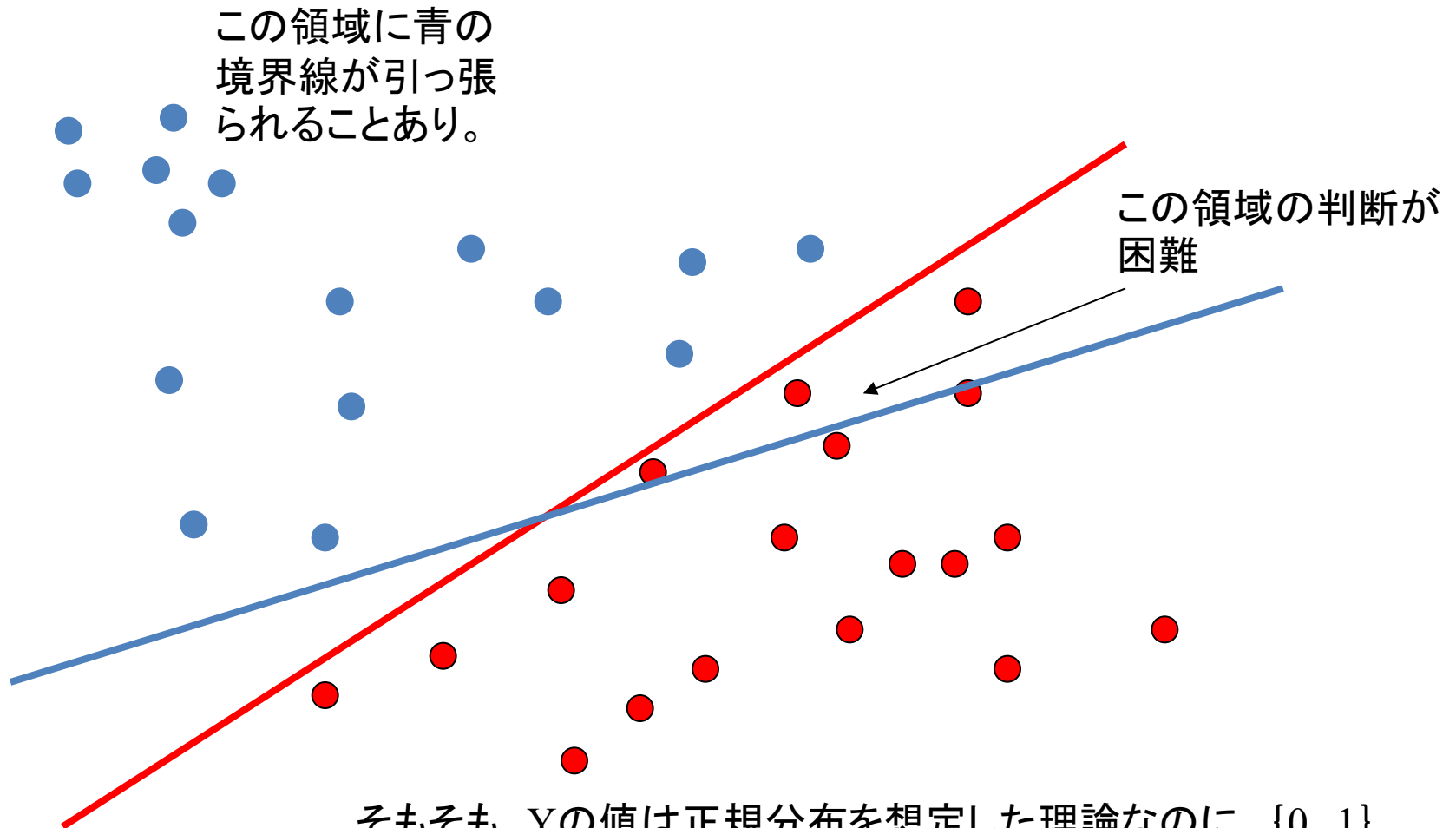
$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = \frac{1}{1 + \exp\left(-(\mathbf{w}^T \mathbf{x} + w_0)\right)}$$

$$\text{where } \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2)$$

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

なので、この分布を教師データから学習できた。

2乗誤差最小化の線形識別の問題点



そもそも、 Y の値は正規分布を想定した理論なのに、 $\{0, 1\}$ の2値しかとらないとして2乗誤差最小化を当てはめたところに無理がある。