

# パターン認識

---

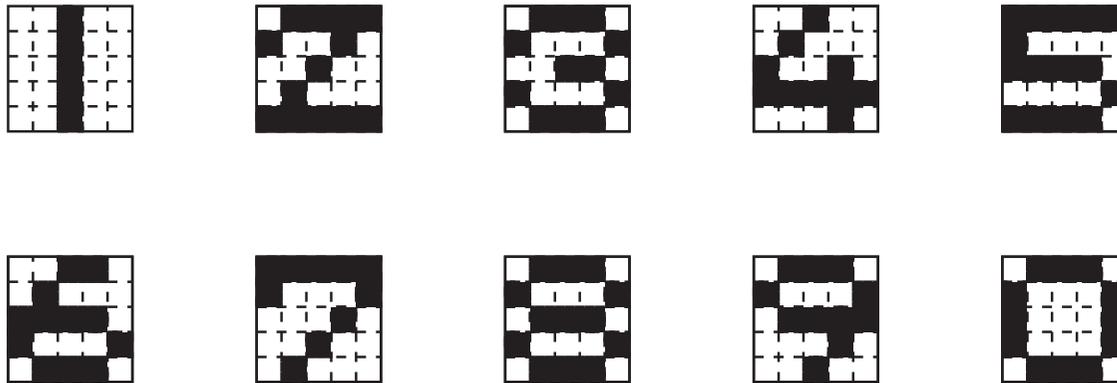
- パターン認識: 「パターンをそれが属すべきカテゴリに対応づける操作」
  - パターン: 「空間的または時間的に観測可能な事象であって, 観測された対象どうしが同一であるか(または似ているか) 否かを判定できるような性質を備えているもの」
  - カテゴリ(またはクラス): 「パターン認識の結果, 同等とみなされるパターンの集合概念」
  - パターンベクトル: パターンを表す  $n$  個の数値の組

$$\boldsymbol{x} = (x_1, x_2, \dots, x_n)$$

- 特徴ベクトル: パターンに関する本質的な情報を含む少数個の数値の組. 位相を崩さずに次元数を極力減らす.

# 特徴ベクトルと特徴空間の例

---



5 × 5メッシュによる2値パターン

- パターン

$$\boldsymbol{x} = (x_1, x_2, \dots, x_{25})$$

- 特徴量

$$\begin{cases} x_j = 1 & (\text{黒: 文字部分}) \\ x_j = 0 & (\text{白: 背景部分}) \end{cases} \quad (1 \leq j \leq 25)$$

- パターンの組み合わせ:  $2^{25} = 33554432$  通り
- 出力カテゴリ: 0 から 9 までの数字

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_{10}\} \quad (c = 10)$$

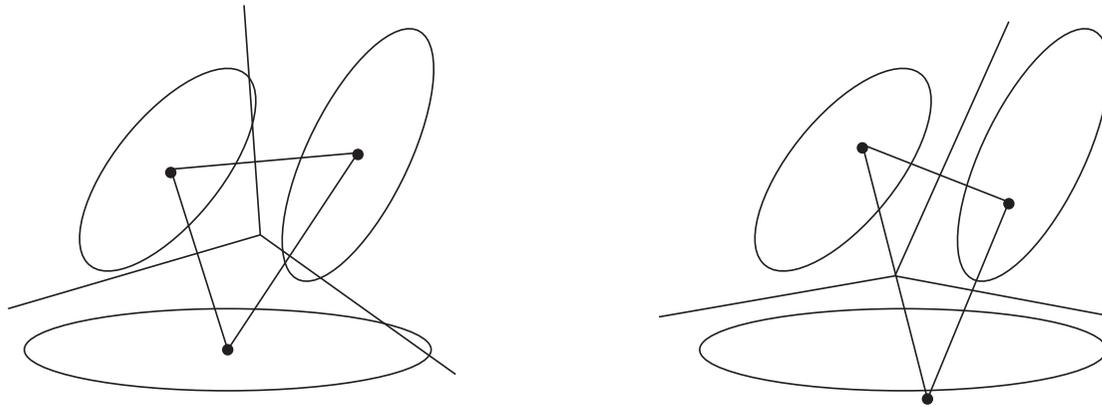
## 決定則の例

---

- 最近傍決定則: 入力パターンともっとも距離の近いプロトタイプの属するカテゴリを選ぶ。

$$\min_{p=1,\dots,n} \{D(\mathbf{x}, \mathbf{x}_p)\} = D(\mathbf{x}, \mathbf{x}_k) \Rightarrow \mathbf{x} \in \omega_k$$

- 問題点
  - 特徴空間の分割法 (3 と 8 は似ている)
  - プロトタイプの選び方 (6 と回転した 9 の識別)



領域分割とプロトタイプ

# 重みベクトル

---

- 2分識別器 (正のクラスと負のクラスを識別):

$$f : X \subset \mathbf{R}^n \rightarrow \mathbf{R} \quad \text{s.t.} \quad \begin{cases} \text{if } \mathbf{x} \in P : f(\mathbf{x}) \geq 0 \\ \text{otherwise} : f(\mathbf{x}) < 0 \end{cases}$$

ただし  $P \subset X$  は正のクラス

- 教師データ:

$$y_i = \begin{cases} 1 & \text{if } \mathbf{x} \in P \\ -1 & \text{otherwise} \end{cases}$$

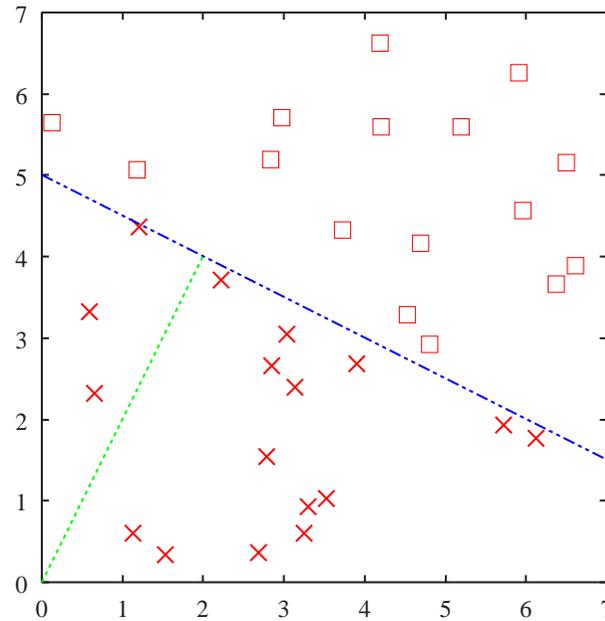
- 線形識別器

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b = \sum_{i=1}^n w_i x_i + b$$

ただし  $(\mathbf{w}, b) \in (\mathbf{R}^n, \mathbf{R})$

- 識別結果:  $y_i f(\mathbf{x}) > 0$  ? (もし正なら正解)
- 線形識別可能性: 線形識別器でカテゴリズ可能かどうか
- 識別面:  $\mathbf{R}^n$  における  $\mathbf{R}^{n-1}$  超平面
  - パラメータ数:  $n + 1$
  - $n$  次元は  $\mathbf{R}^n$  における平面の法線ベクトル  $\mathbf{w}$  の次元
  - 1 次元は平面の原点からの距離  $b$

# 線形識別器

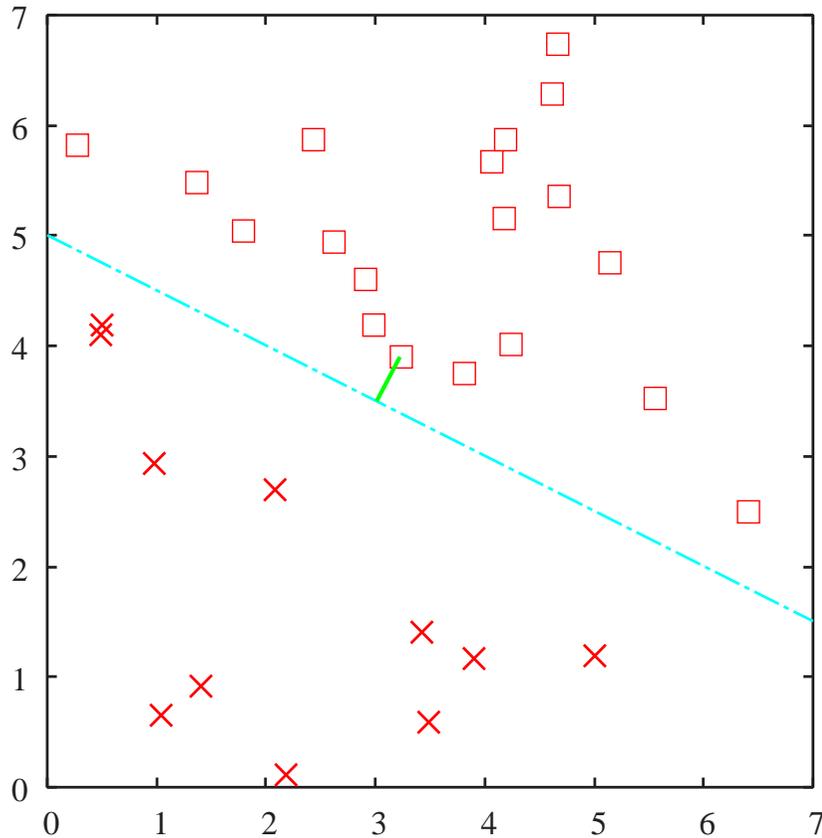


- 識別式 ( $\|\mathbf{w}\| = 1$ )

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b = \sum_{i=1}^n w_i x_i + b$$

- 識別面の方程式:  $f(\mathbf{x}) = 0$
- 原点から識別面までの距離:  $|b|$
- 識別面の法線ベクトル:  $\mathbf{w}$
- 識別面から点  $\mathbf{x}_j$  までの (符合つき) 距離:  $f(\mathbf{x}_j)$

# マージン



- 超平面  $(w, b)$  に関するトレーニングデータ  $(x_i, y_i)$  のマージン (識別面からの距離)

$$\gamma_i = y_i(\langle w, x_i \rangle + b)$$

- $\gamma_i > 0$  ならば  $(x_i, y_i)$  に関して正しい識別
- 超平面  $(w, b)$  に関するトレーニングセット  $S$  のマージン: トレーニングデータに関するマージンの最小値
- トレーニングセット  $S$  のマージン: すべての超平面  $(w, b)$  に関するトレーニングセット  $S$  のマージンの最小値

# パーセプトロン

---

- F. Rosenblatt: The perceptron: a probabilistic model for information storage and organization in the brain, *Psychological Review*, 65: 386-408, 1959 1960.
- M. L. Minsky and S. A. Papert: Perceptrons, MIT Press, 1969. Expanded Edition 1990.

- 入力空間:  $X \subset \mathbf{R}^n$
- 出力カテゴリ:  $Y = -1, 1$
- トレーニングデータ:  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)) \subset (X \times Y)^\ell$
- 判別:  $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 0$  ?
- 重み更新:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta y_i \mathbf{x}_i$$

$$b_{k+1} = b_k + \eta y_i R^2$$

$$R = \max_{1 \leq i \leq \ell} \|\mathbf{x}_i\|$$

# アルゴリズム

---

Given linearly separable training set  $S$  and learning rate  $\eta > 0$

$$\mathbf{w}_0 = \mathbf{0}; b_0 = 0; k = 0$$

$$R = \max_{1 \leq i \leq \ell} \|\mathbf{x}_i\|$$

do

  for  $i = 1$  to  $\ell$

    if  $y_i(\langle \mathbf{w}_k, \mathbf{x}_i \rangle + b_k) \leq 0$  then

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta y_i \mathbf{x}_i$$

$$b_{k+1} = b_k + \eta y_i R^2$$

$$k = k + 1$$

    end if

  end for

until no mistakes made within the for loop

return  $k, (\mathbf{w}_k, b_k)$

where  $k$  is the number of mistakes

## 収束定理

---

(Novikoff) Let  $S$  be a non-trivial training set, and let

$$R = \max_{1 \leq i \leq \ell} \|\mathbf{x}_i\|.$$

Suppose that there exists a vector  $\mathbf{w}_{opt}$  such that  $\|\mathbf{w}_{opt}\| = 1$  and

$$y_i(\langle \mathbf{w}_{opt}, \mathbf{x}_i \rangle + b_{opt}) \geq \gamma$$

for  $1 \leq i \leq \ell$ . Then the number of mistakes made by the on-line perceptron algorithm on  $S$  is at most

$$\left(\frac{2R}{\gamma}\right)^2.$$

- 線形識別可能性を仮定すると間違いの回数の上限が与えられる
- 線形識別可能でない場合は修正を繰り返し収束しない

## 証明

---

Define  $\hat{\mathbf{x}}_i = (\mathbf{x}_i^T, R)^T$ ,  $\hat{\mathbf{w}} = (\mathbf{w}^T, b/R)^T$ . Let  $\hat{\mathbf{w}}_{t-1}$  be the augmented weight vector prior to the  $t$ th mistake. The  $t$ th update is performed when

$$y_i \langle \hat{\mathbf{w}}_{t-1}, \hat{\mathbf{x}}_i \rangle = y_i (\langle \mathbf{w}_{t-1}, \mathbf{x}_i \rangle + b_{t-1}) \leq 0$$

where  $(\mathbf{x}_i, y_i)$  is the point incorrectly classified by  $\hat{\mathbf{w}}_{t-1}$ . The update is the following:

$$\hat{\mathbf{w}}_t = \begin{bmatrix} \mathbf{w}_t \\ b_t/R \end{bmatrix} = \begin{bmatrix} \mathbf{w}_{t-1} \\ b_{t-1}/R \end{bmatrix} + \eta y_i \begin{bmatrix} \mathbf{x}_i \\ R \end{bmatrix} = \hat{\mathbf{w}}_{t-1} + \eta y_i \hat{\mathbf{x}}_i$$

where  $b_t = b_{t-1} + \eta y_i R^2$ . Since the margin is  $\gamma$ , we have

$$\langle \hat{\mathbf{w}}_t, \hat{\mathbf{w}}_{opt} \rangle = \langle \hat{\mathbf{w}}_{t-1}, \hat{\mathbf{w}}_{opt} \rangle + \eta y_i \langle \hat{\mathbf{x}}_i, \hat{\mathbf{w}}_{opt} \rangle \geq \langle \hat{\mathbf{w}}_{t-1}, \hat{\mathbf{w}}_{opt} \rangle + \eta \gamma$$

and this implies that

$$\langle \hat{\mathbf{w}}_t, \hat{\mathbf{w}}_{opt} \rangle \geq t \eta \gamma$$

Similarly, we have

$$\begin{aligned} \|\hat{\mathbf{w}}_t\|^2 &= \|\hat{\mathbf{w}}_{t-1}\|^2 + 2\eta y_i \langle \hat{\mathbf{w}}_{t-1}, \hat{\mathbf{x}}_i \rangle + \eta^2 \|\hat{\mathbf{x}}_i\|^2 \\ &\leq \|\hat{\mathbf{w}}_{t-1}\|^2 + \eta^2 \|\hat{\mathbf{x}}_i\|^2 \leq \|\hat{\mathbf{w}}_{t-1}\|^2 + \eta^2 (\|\mathbf{x}_i\|^2 + R^2), \end{aligned}$$

which implies  $\|\hat{\mathbf{w}}_t\|^2 \leq \|\hat{\mathbf{w}}_{t-1}\|^2 + 2\eta^2 R^2$ . Thus we have

$$\|\hat{\mathbf{w}}_t\|^2 \leq 2t\eta^2 R^2.$$

## 証明(つづき)

---

The two inequalities combined give the ‘squeezing’ relations

$$\|\mathbf{w}_{opt}\| \sqrt{2t\eta R} \geq \|\mathbf{w}_{opt}\| \|\mathbf{w}_t\| \geq \langle \hat{\mathbf{w}}_t, \hat{\mathbf{w}}_{opt} \rangle \geq t\eta\gamma,$$

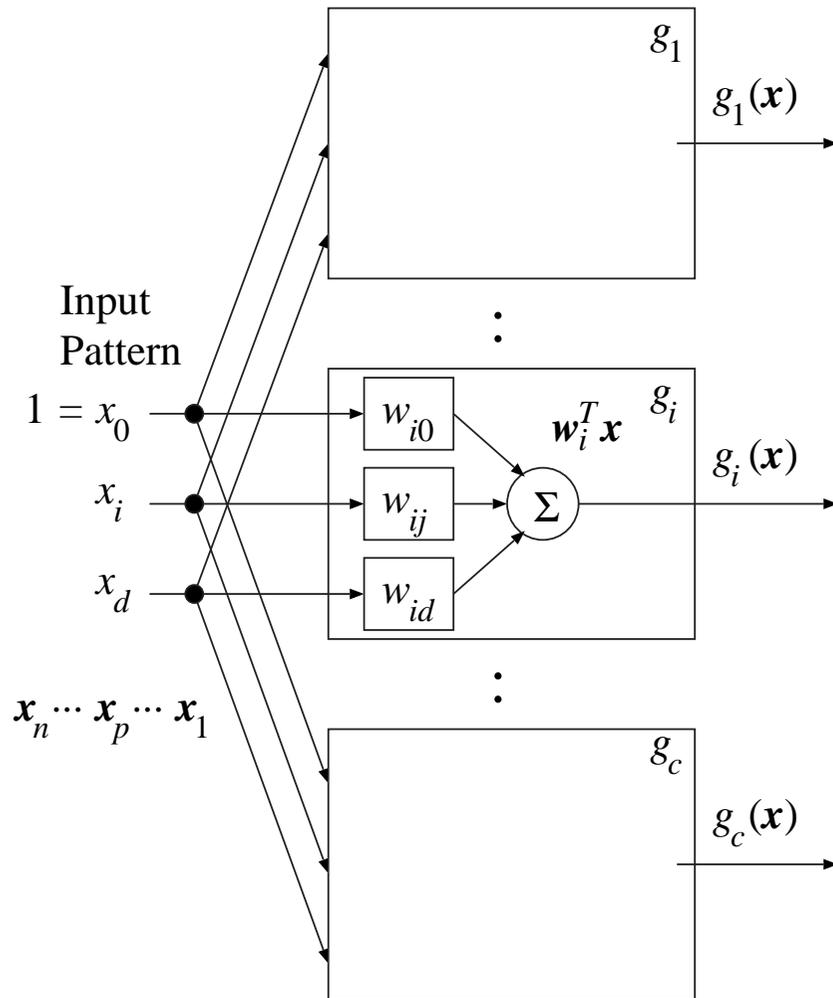
which together imply the bound

$$t \leq 2 \left( \frac{R}{\gamma} \right)^2 \|\hat{\mathbf{w}}_{opt}\|^2 \leq \left( \frac{2R}{\gamma} \right)^2$$

since  $b_{opt} \leq R$  and thus  $\|\hat{\mathbf{w}}_{opt}\|^2 \leq \|\mathbf{w}_{opt}\|^2 + 1 = 2$ .

- 線形識別可能性を仮定すれば間違いの回数に上限が存在
- つまり有限回の繰り返しで収束
- 一方，線形識別可能でなければ修正を繰り返し収束しない
- したがって，線形識別可能性の判断には使えない
- 途中で打ち切った場合，そのときの重みが最適かどうかわからない
- そもそも最適性の評価基準は？

# 誤差評価に基づく学習



- カテゴリ:  $\omega_1, \dots, \omega_c$
- 学習パターン:  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
- 教師ベクトル:  $\mathbf{b}_1, \dots, \mathbf{b}_n$

$$\mathbf{b}_p = [b_{1p}, \dots, b_{cp}]^T$$

- 出力ベクトル:

$$\mathbf{g}(\mathbf{x}_p) = [g_1(\mathbf{x}_p), \dots, g_c(\mathbf{x}_p)]^T$$

- 線形識別関数

$$g_i(\mathbf{x}_p) = w_{i0} + \sum_{j=1}^d w_{ij} x_{pj} = \mathbf{w}_i^T \mathbf{x}_p$$

$$\mathbf{w}_i = [w_{i0}, w_{i1}, \dots, w_{id}]^T,$$

$$\mathbf{x} = [1, x_{p1}, \dots, x_{pd}]^T$$

## 評価関数

---

- 誤差:  $\boldsymbol{\varepsilon}_p = [\varepsilon_{1p}, \dots, \varepsilon_{cp}]^T = \mathbf{g}(\mathbf{x}_p) - \mathbf{b}_p$

$$\varepsilon_{ip} = g_i(\mathbf{x}_p) - b_{ip}$$

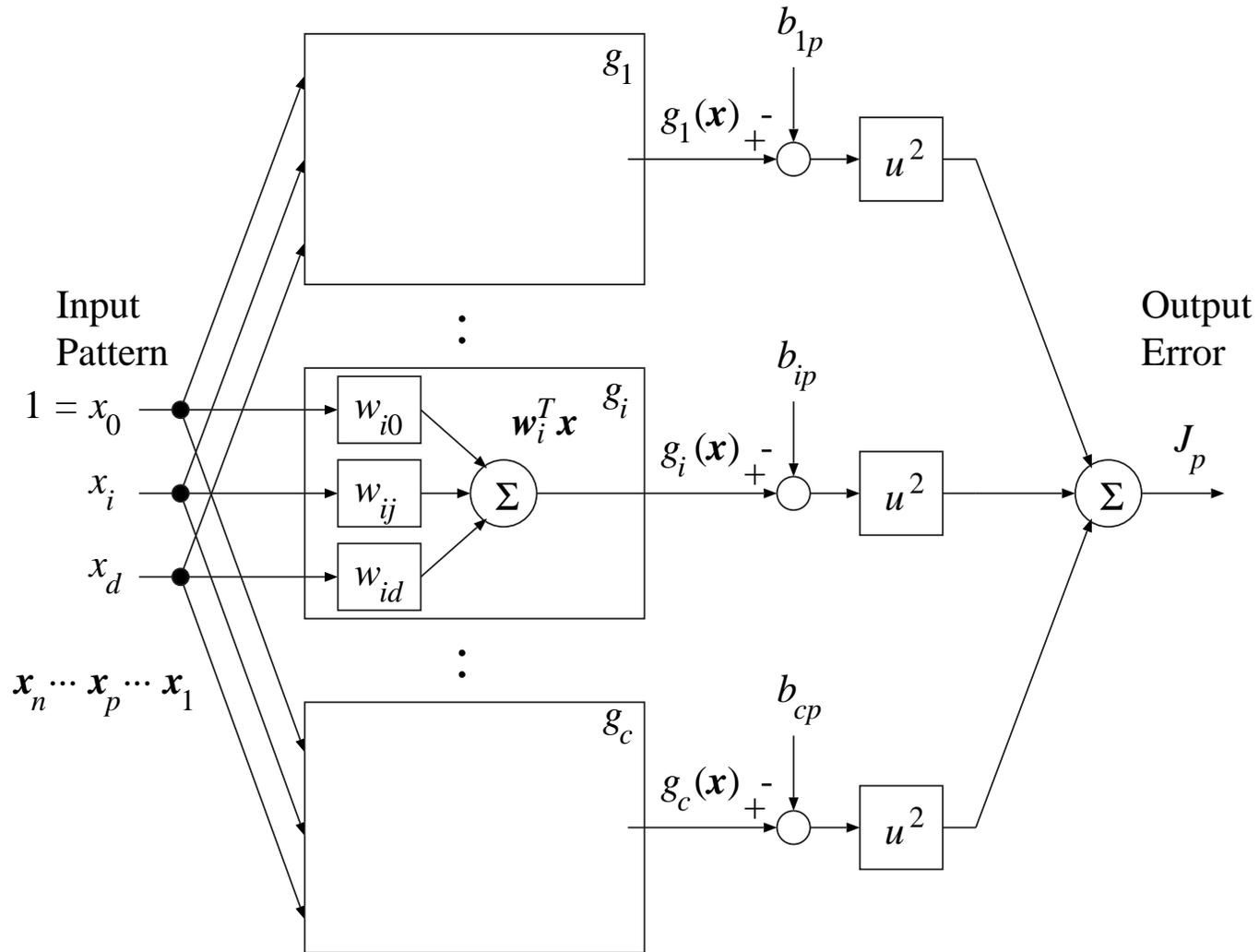
- パターン  $\mathbf{x}_p$  に対する評価関数:  $J_p$

$$J_p(\mathbf{w}_1, \dots, \mathbf{w}_c) = \frac{1}{2} \|\boldsymbol{\varepsilon}_p\|^2 = \frac{1}{2} \sum_{i=1}^c \varepsilon_{ip}^2 = \frac{1}{2} \sum_{i=1}^c (\mathbf{w}_i^T \mathbf{x}_p - b_{ip})^2$$

- 全パターンに対する評価関数:  $J$

$$J(\mathbf{w}_1, \dots, \mathbf{w}_c) = \sum_{p=1}^n J_p(\mathbf{w}_1, \dots, \mathbf{w}_c) = \frac{1}{2} \sum_{p=1}^n \sum_{i=1}^c (\mathbf{w}_i^T \mathbf{x}_p - b_{ip})^2$$

# 評価機構をもつ識別器



## 最適重み

---

- 最適重み: 評価関数  $J$  を最小にする  $w_i$  ( $i = 1, \dots, c$ )

$$\frac{\partial J}{\partial w_i} = 0 \quad (i = 1, \dots, c)$$

- 全パターンに対する評価関数:  $J$

$$J(\mathbf{w}_1, \dots, \mathbf{w}_c) = \frac{1}{2} \sum_{p=1}^n \sum_{i=1}^c (\mathbf{w}_i^T \mathbf{x}_p - b_{ip})^2$$

- 最適重み  $w_i$  は次式をみたす

$$\begin{aligned} 0 &= \frac{\partial J}{\partial w_i} = \sum_{p=1}^n (\mathbf{w}_i^T \mathbf{x}_p - b_{ip}) \mathbf{x}_p = [\mathbf{x}_1, \dots, \mathbf{x}_n] \begin{bmatrix} \mathbf{w}_i^T \mathbf{x}_1 - b_{i1} \\ \vdots \\ \mathbf{w}_i^T \mathbf{x}_n - b_{in} \end{bmatrix} \\ &= [\mathbf{x}_1, \dots, \mathbf{x}_n] \left( \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \mathbf{w}_i - \begin{bmatrix} b_{i1} \\ \vdots \\ b_{in} \end{bmatrix} \right) = X^T (X \mathbf{w}_i - \bar{\mathbf{b}}_i) \end{aligned}$$

ただし ( $X$  はパターン行列と呼ばれる)

$$X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T, \quad \bar{\mathbf{b}}_i = [b_{i1}, \dots, b_{in}]^T$$

## 最適重み

---

- 注意:  $\bar{b}_i$  は教師ベクトル  $b_p$  とは以下の関係がある .

$$[\mathbf{b}_1, \dots, \mathbf{b}_n] = \begin{bmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{c1} & \cdots & b_{cn} \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{b}}_1^T \\ \vdots \\ \bar{\mathbf{b}}_c^T \end{bmatrix}$$

- 最適重み

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

## 逐次法

---

- 最急降下法 (パターンが提示されるたびに修正を行う)

$$\mathbf{w}'_i = \mathbf{w}_i - \rho \frac{\partial J_p}{\partial \mathbf{w}_i} \quad (i = 1, \dots, c)$$

$$J_p(\mathbf{w}_1, \dots, \mathbf{w}_c) = \frac{1}{2} \sum_{i=1}^c (\mathbf{w}_i^T \mathbf{x}_p - b_{ip})^2$$

ここで  $\varepsilon_{ip} = \mathbf{w}_i^T \mathbf{x}_p - b_{ip}$  に注意すると

$$\frac{\partial J_p}{\partial \mathbf{w}_i} = \frac{\partial J_p}{\partial \varepsilon_{ip}} \frac{\partial \varepsilon_{ip}}{\partial \mathbf{w}_i} = \varepsilon_{ip} \mathbf{x}_p$$

であるから

$$\mathbf{w}'_i = \mathbf{w}_i - \rho \varepsilon_{ip} \mathbf{x}_p = \mathbf{w}_i - \rho (\mathbf{w}_i^T \mathbf{x}_p - b_{ip}) \mathbf{x}_p \quad (i = 1, \dots, c)$$

を得る . これを Widrow-Hoff の学習則という .

## パーセプトロン(再)

---

- 学習:  $\mathcal{X}_i$  に属するすべての  $x$  に対して

$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad (j = 1, \dots, c, j \neq i)$$

が成り立つように重み  $w_i$  を決定すること .

- 学習アルゴリズム

1. 重みベクトル  $w_i$  の初期値を適当に設定する .
2.  $\mathcal{X}_i$  のなかから学習パターンをひとつ選ぶ .
3. 識別関数

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x}$$

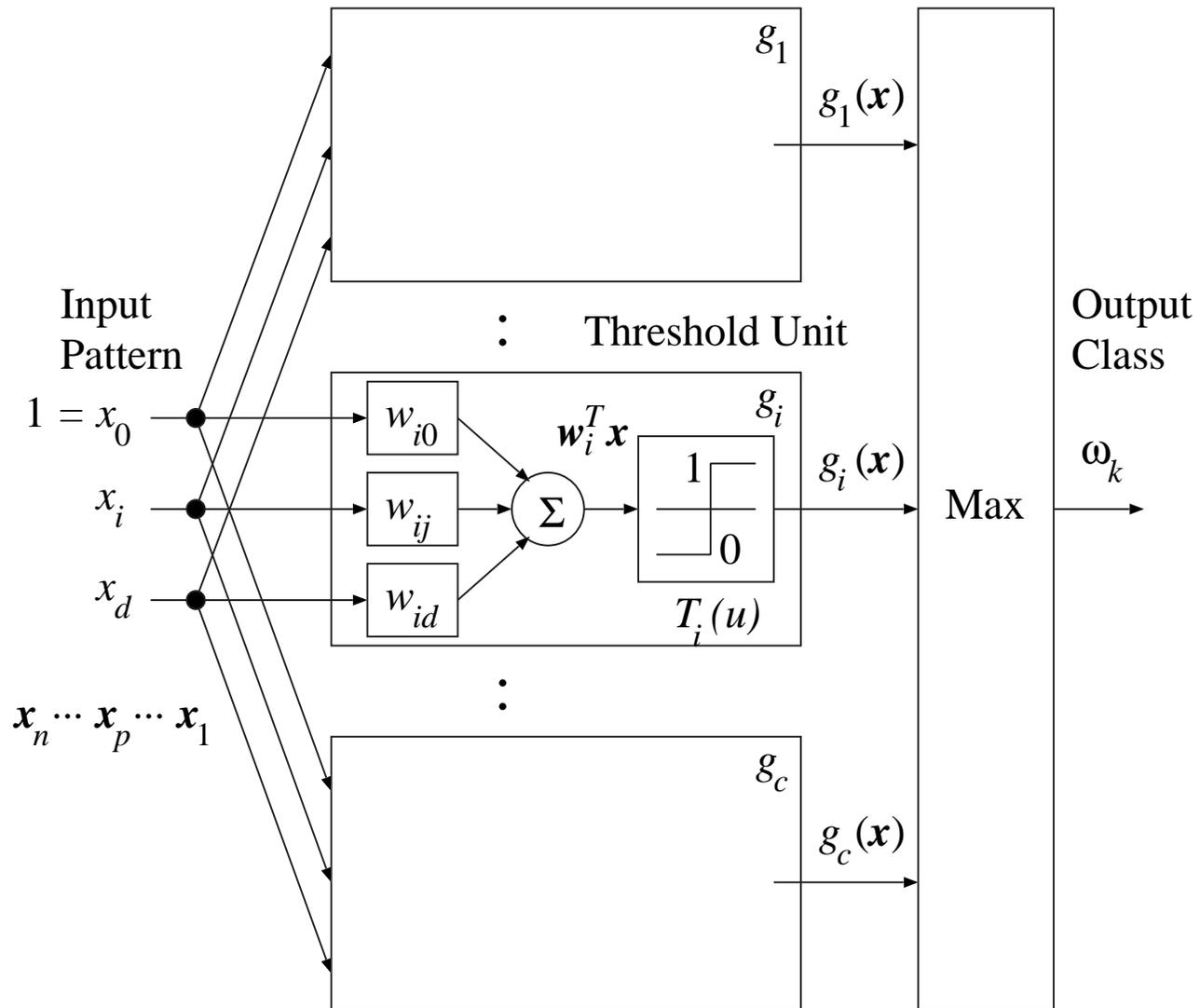
によって識別を行い, 誤識別を生じた場合 ( $\omega_i$  に属するパターンを  $\omega_j$  と誤ったとする) のみ次式により重みベクトルを修正する .

$$\mathbf{w}'_i = \mathbf{w}_i + \rho \mathbf{x}$$

$$\mathbf{w}'_j = \mathbf{w}_j - \rho \mathbf{x}$$

4. 上の処理, 2, 3 を  $\mathcal{X}_i$  の全パターンに対して繰り返す .
5.  $\mathcal{X}_i$  の全パターンを正しく識別できたら終了, 誤りがあるときには 2 に戻る .

# パーセプトロン(ブロック図)



## 誤差評価とパーセプトロン

---

- 教師信号

$$b_{ip} = \begin{cases} 1 & (\mathbf{x}_p \in \omega_i) \\ 0 & (\mathbf{x}_p \notin \omega_i) \end{cases} \quad (i = 1, \dots, c)$$

- 出力にしきい値関数  $T_i$  を施す

$$T_i(u) = \begin{cases} 1 & (u > 0) \\ 0 & (u \leq 0) \end{cases} \quad (i = 1, \dots, c)$$

- Widrow-Hoff 学習規則 (パターン  $\mathbf{x}_p \in \omega_i$  を  $\omega_j$  と誤認識したとき)

$$\begin{cases} \mathbf{w}'_i = \mathbf{w}_i + \rho(g_i(\mathbf{x}_p) - b_{ip})\mathbf{x}_p = \mathbf{w}_i + \rho\mathbf{x}_p \\ \mathbf{w}'_j = \mathbf{w}_j + \rho(g_j(\mathbf{x}_p) - b_{jp})\mathbf{x}_p = \mathbf{w}_j - \rho\mathbf{x}_p \end{cases}$$

- パーセプトロンは Widrow-Hoff 学習規則の一種 !?