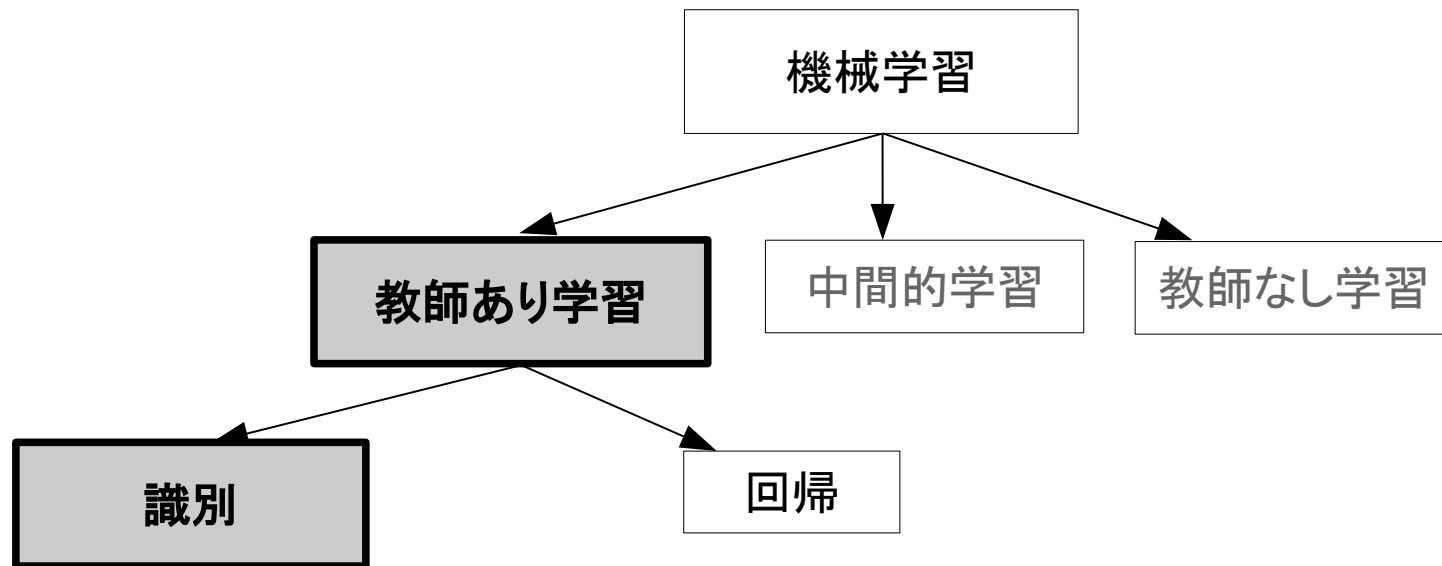
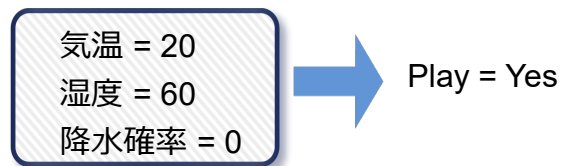


5. 識別 — 生成モデルと識別モデル—

- 問題設定
 - 教師あり学習
 - 数値入力 → カテゴリ出力



- 数値特徴



5.1 数値特徴に対する「教師あり・識別」問題の定義

- 教師あり学習のデータ

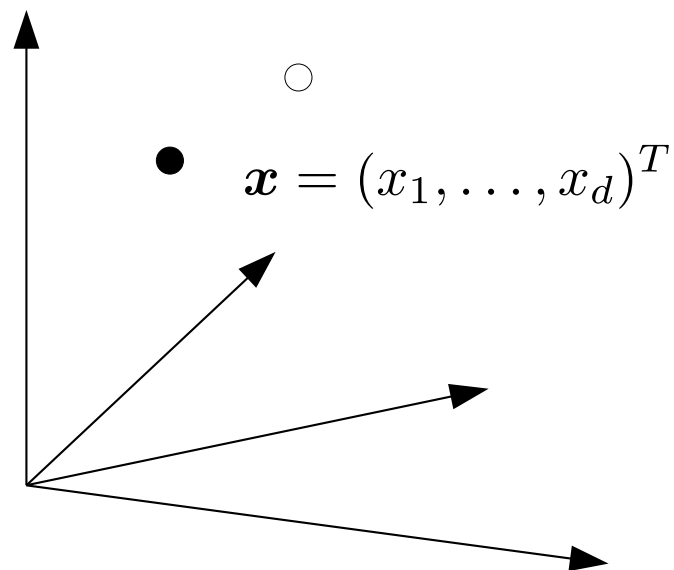
- 特徴ベクトル \mathbf{x} と正解情報 y のペア

$$\{(\mathbf{x}_i, y_i)\}, \quad i = 1 \dots N$$

- \mathbf{x} は次元数 d の固定長ベクトル、 y はカテゴリ

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

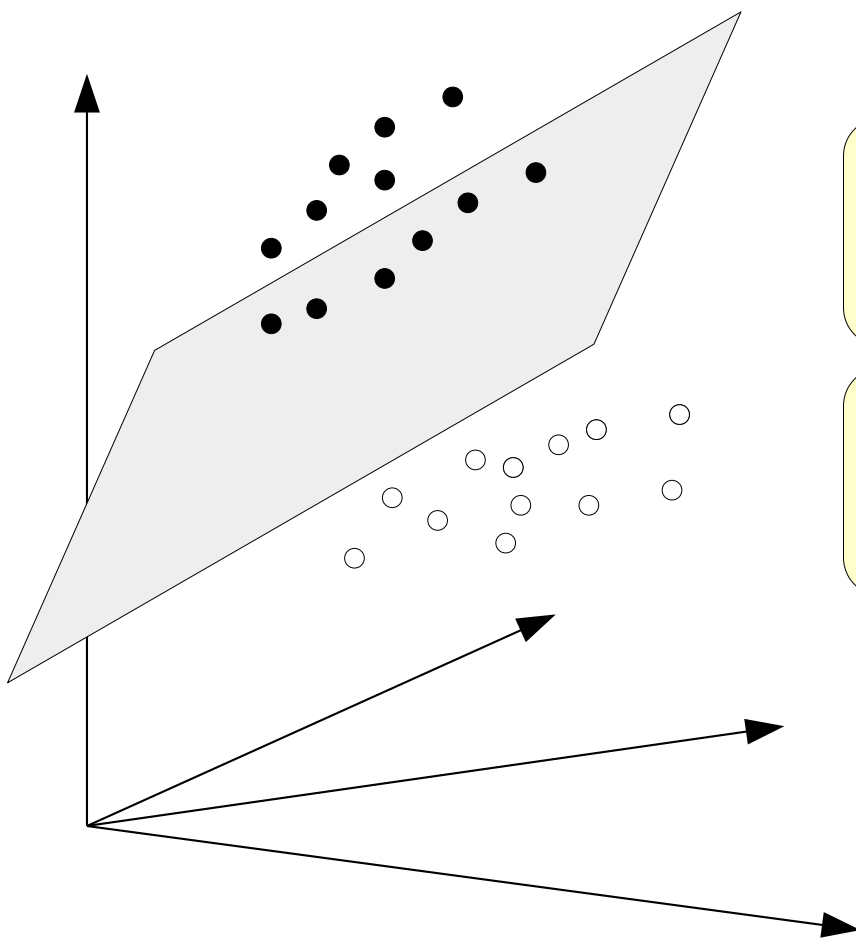
- \mathbf{x} は d 次元空間（特徴空間）上の点と見なせる



- $y = \text{positive}$ (正例)
- $y = \text{negative}$ (負例)

5.1 数値特徴に対する「教師あり・識別」問題の定義

- 数値特徴に対する識別問題＝境界面の設定
 - 各クラスの確率分布を求めることで、結果として境界面（等確率となる点の集合）が定まる場合も含む



クラスが比較的きれいに分離している
→ 小数のパラメータで各クラスを表現可能
⇒ 統計的手法（第5章）

クラス境界が複雑
高次元へマッピング⇒ **SVM（第7章）**
非線形識別面 ⇒ **ニューラルネット（第8章）**

4.1 統計的識別とは

復習

- 最大事後確率則による識別

$$C_{MAP} = \arg \max_i P(\omega_i | \mathbf{x})$$

\mathbf{x} : 特徴ベクトル
 ω_i ($1 \leq i \leq c$) : クラス

- データから直接的にこの確率を求めるのは難しい

- ベイズの定理 $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

$$C_{MAP} = \arg \max_i P(\omega_i | \mathbf{x})$$

$$= \arg \max_i \frac{P(\mathbf{x} | \omega_i) P(\omega_i)}{P(\mathbf{x})}$$

ベイズの定理

$$= \arg \max_i \underbrace{P(\mathbf{x} | \omega_i)}_{\text{尤度}} \underbrace{P(\omega_i)}_{\text{事前確率}}$$

上式の分母は
判定に寄与しない

尤度

事前確率

4.1 統計的識別とは

復習

- 事前確率 $P(\omega_i)$
 - 特徴ベクトル \mathbf{x} を観測する前の、各クラスの起こりやすさ
- 事前確率の最尤推定

$$P(\omega_i) = \frac{n_i}{N}$$

N : 全データ数、 n_i : クラス ω_i のデータ数

4.2.2 ナイーブベイズ識別

復習

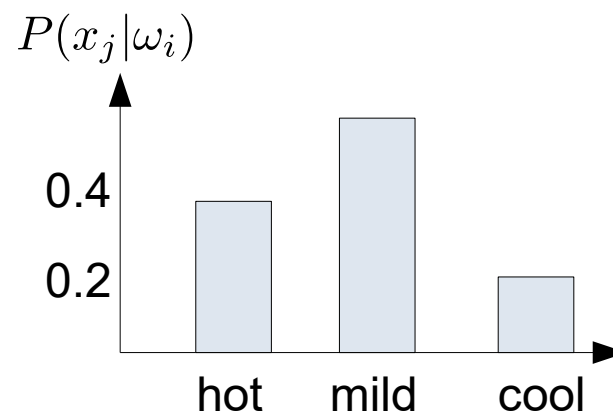
- 尤度計算におけるナイーブベイズの近似
 - すべての特徴が独立であると仮定

$$P(\mathbf{x}|\omega_i) = P(x_1, \dots, x_d|\omega_i)$$

$$\approx \prod_{j=1}^d P(x_j|\omega_i)$$

$$C_{NB} = \arg \max_i P(\omega_i) \prod_{j=1}^d P(x_j|\omega_i)$$

- x_j がカテゴリ特徴のとき
 - $y = \omega_i$ のデータから、カテゴリ分布を最尤推定



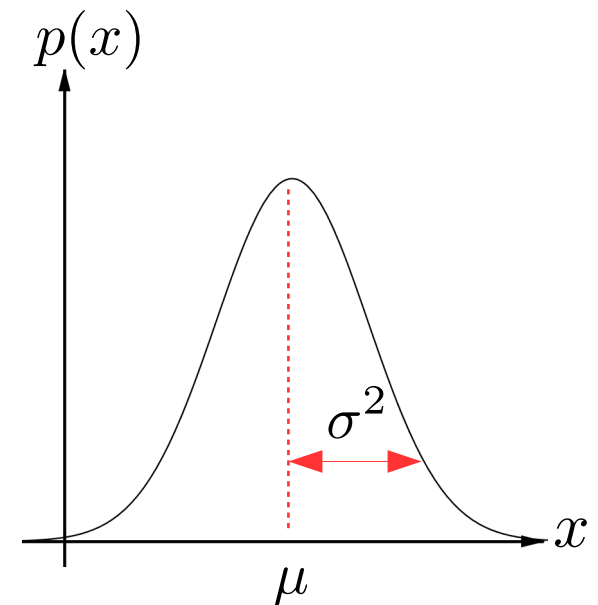
5.2 生成モデル

5.2.1 数値特徴に対するナイーブベイズ識別

- 確率密度関数 $p(x_j | \omega_i)$ の推定
 - 関数をクラス毎に求めるので ω_i は省略
 - 関数の形は正規分布を仮定

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

- データの対数尤度を最大とする
平均 μ と分散 σ^2 を求める



正規分布とは

- 離散型二項分布の例

- n 枚のコインを投げた時の、表の枚数の度数

- $n=1$ 1 1

- $n=2$ 1 2 1

- $n=3$ 1 3 3 1

- $n=4$ 1 4 6 4 1

- $n=5$ 1 5 10 10 5 1

- ...

- $n \rightarrow \infty$ の時が正規分布

- さまざまな要因が重なって生じた数値がこの分布になる

- 例) テストの点数 (n 問からなる問題の結果)

- 身長 (遺伝、食事、運動、睡眠、 ...)

5.2.1 数値特徴に対するナイーブベイズ識別

- データの対数尤度（最大化したい）

$$\mathcal{L}(D) = \log P(D|\mu, \sigma^2) = \sum_{i=1}^N \log P(x_i|\mu, \sigma^2)$$

- 正規分布の式を当てはめる

$$\mathcal{L}(D) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

- μ で偏微分して 0 とおく

$$\frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

求める分布の平均はデータの平均、分散はデータの分散
というごく当たり前の結果

- σ^2 で偏微分して 0 とおく

$$-\frac{N}{2} \frac{1}{\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^N (x_i - \mu)^2 = 0 \Rightarrow \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

5.2.2 生成モデルの考え方

- 事後確率を求めるにあたって、同時確率を求めている
 - データが生成される様子をモデル化しているとも見ること出来る
 - 事前確率に基づいてクラスを選ぶ
 - そのもとで、特徴ベクトルを出力する

$$\begin{aligned} P(\omega_i|\mathbf{x}) &= \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})} \\ &= \frac{p(\omega_i, \mathbf{x})}{p(\mathbf{x})} \end{aligned}$$

事後確率を求めるより、
難しい問題を解いている
のではないかな？

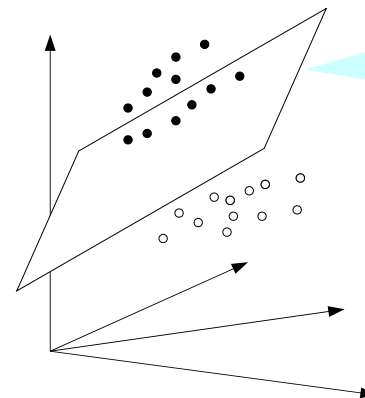
5.3 識別モデル

- 識別関数法

- 確率の枠組みにはとらわれず、 $f_P(\mathbf{x}) > f_N(\mathbf{x})$ ならば \mathbf{x} を positive と判定する関数を推定する
- 2 クラス問題なら $f(\mathbf{x}) = f_P(\mathbf{x}) - f_N(\mathbf{x})$ の正負で判定すればよい ($f(\mathbf{x}) = 0$ が識別面)
- 単層パーセプトロン
 - 識別関数として 1 次式 (= 直線・平面) を仮定

最も単純な識別関数法の実現

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



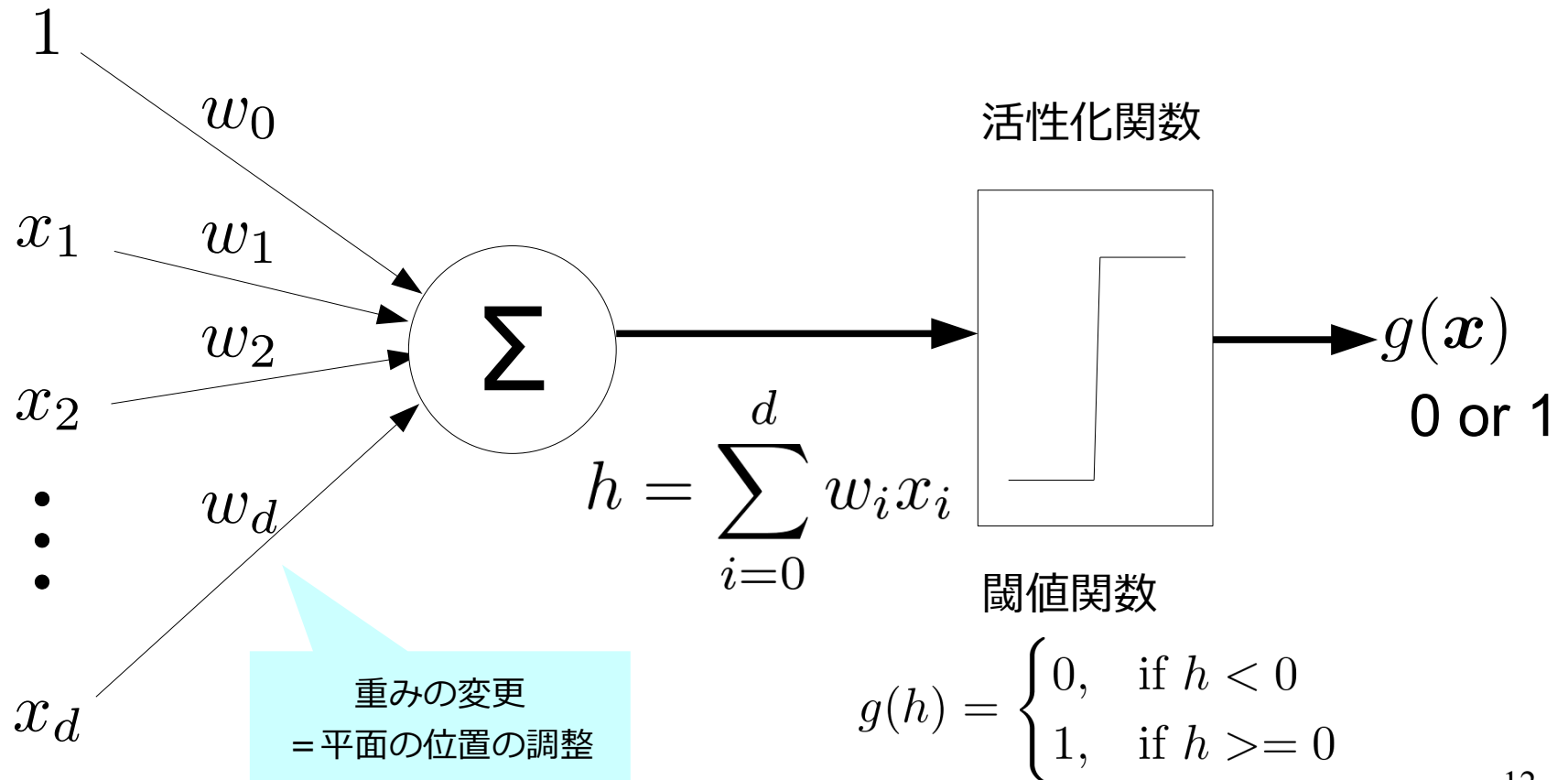
この平面を
求めている
ことになる

5.3.1 誤り訂正学習

- 単層パーセプトロンの定義

以後、 w は w_0 を含む

- $w^T x = 0$ という特徴空間上の超平面を表現



5.3.1 誤り訂正学習

- パーセプトロンの学習規則

1. 重み w の初期値を適当に決める
2. 学習データからひとつ x を選び、 $g(x)$ を計算
3. 誤識別が起きたときのみ、 w を修正する

$$w' = w + \rho x \quad (\text{positive のデータを negative と誤ったとき})$$

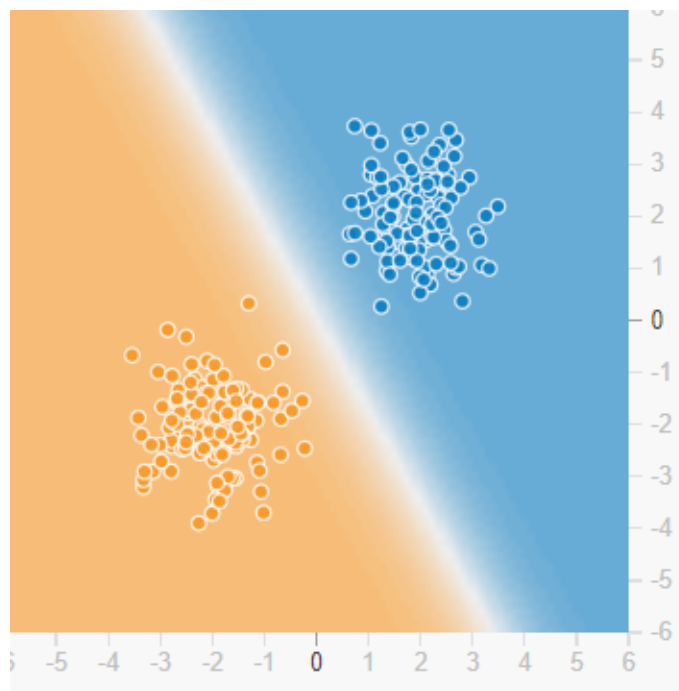
$$w' = w - \rho x \quad (\text{negative のデータを positive と誤ったとき})$$

学習係数

4. 2,3 をすべての学習データについて繰り返す
5. すべて正しく識別できたら終了。そうでなければ 2 へ

5.3.1 誤り訂正学習

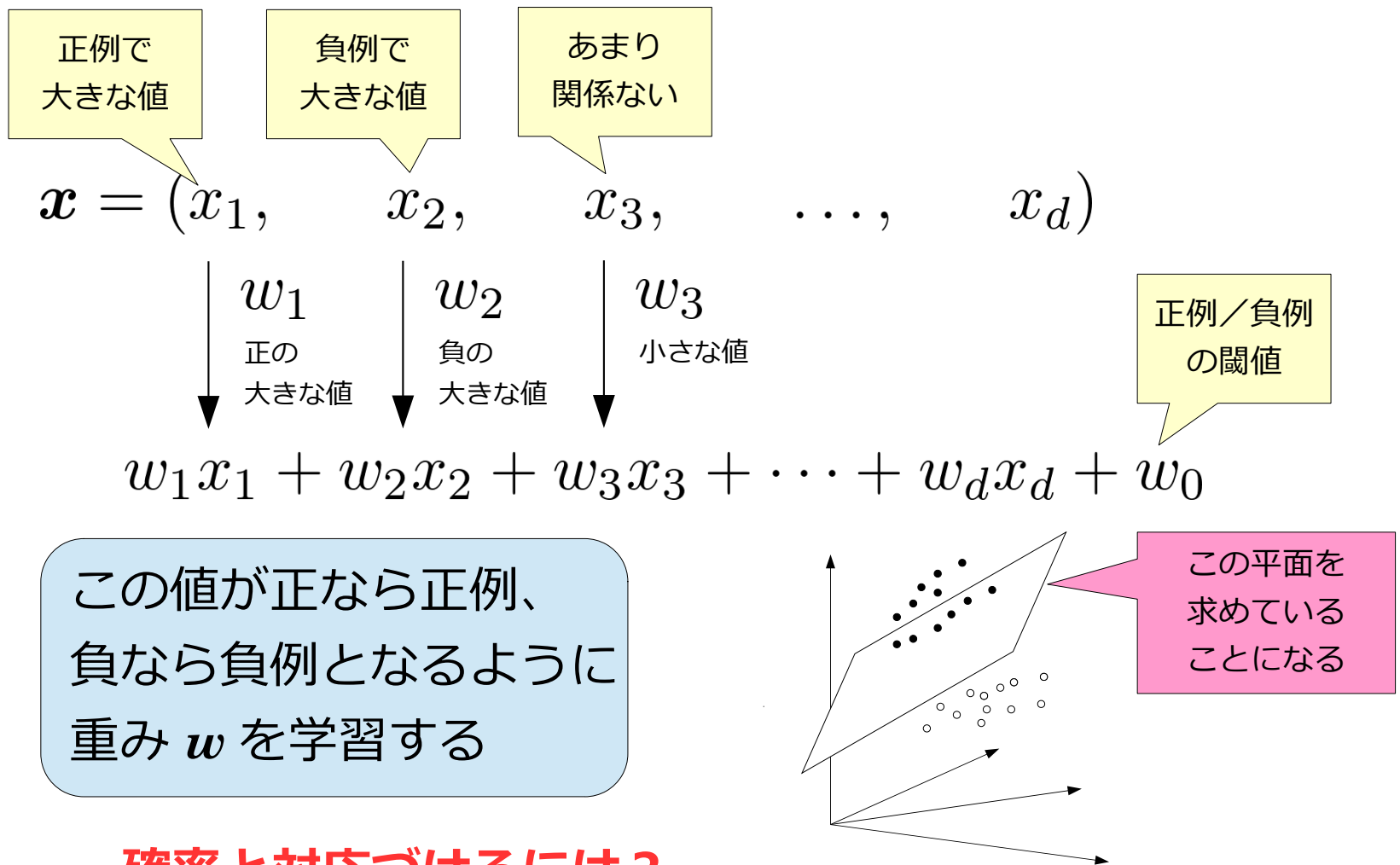
- パーセプトロンの学習規則の適用範囲
 - データが線形分離可能な場合はパーセプトロンの学習規則で学習可能



- 線形分離不可能な場合は終了しない

5.3.3 識別モデルの考え方

- 事後確率を直接求める

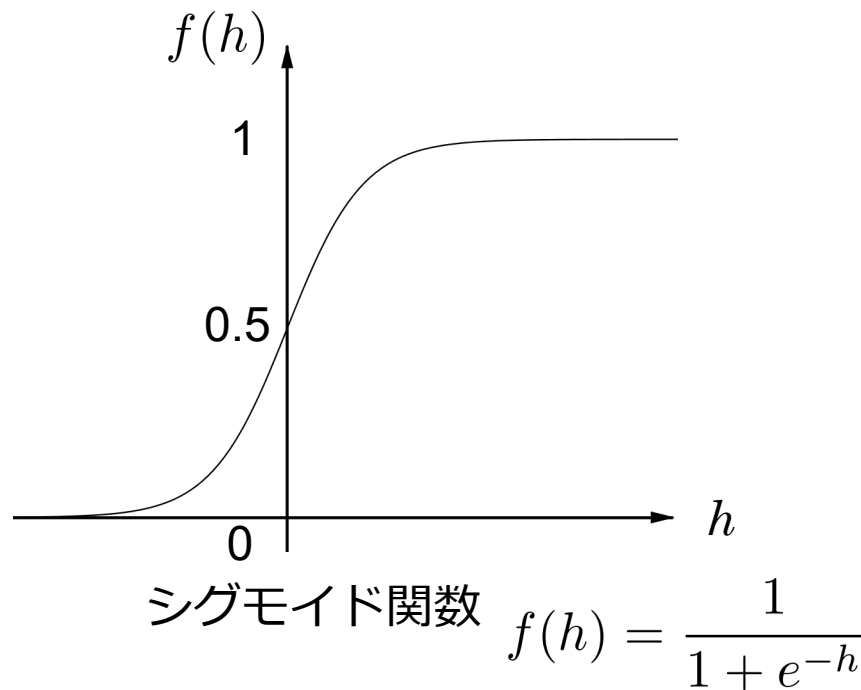


確率と対応づけるには？

5.3.4 ロジスティック識別

- ロジスティック識別
 - 入力が正例である確率

$$P(\text{Positive} \mid \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

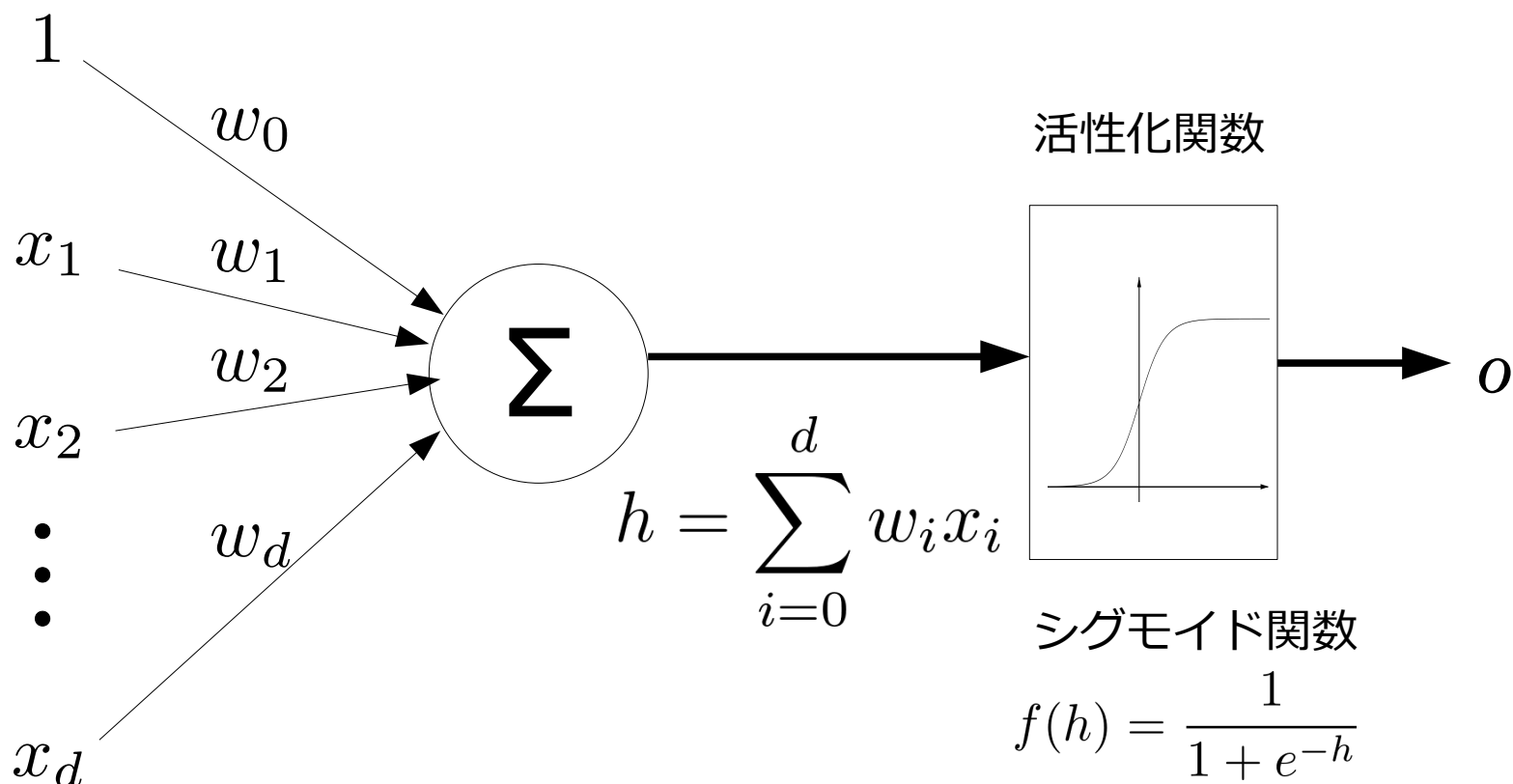


シグモイド関数の利点

- $-\infty \sim +\infty$ の値域を持つものを、順序を変えずに $0 \sim 1$ にマッピングできる
- 微分形が簡単な式になる
 $f'(h) = f(h)(1-f(h))$

5.3.4 ロジスティック識別

- ロジスティック識別の計算ユニット
 - $w^T x$ の値を計算後、シグモイド関数で非線形変換



5.3.4 ロジスティック識別

- 最適化対象：モデルの対数尤度の反数（最小化）

$$E(\boldsymbol{w}) = -\log P(D|\boldsymbol{w})$$

$$= -\log \prod_{\boldsymbol{x}_i \in D} o_i^{y_i} (1 - o_i)^{(1-y_i)}$$

$$o = P(\text{positive} \mid \boldsymbol{x})$$

$y = 0 \text{ or } 1$

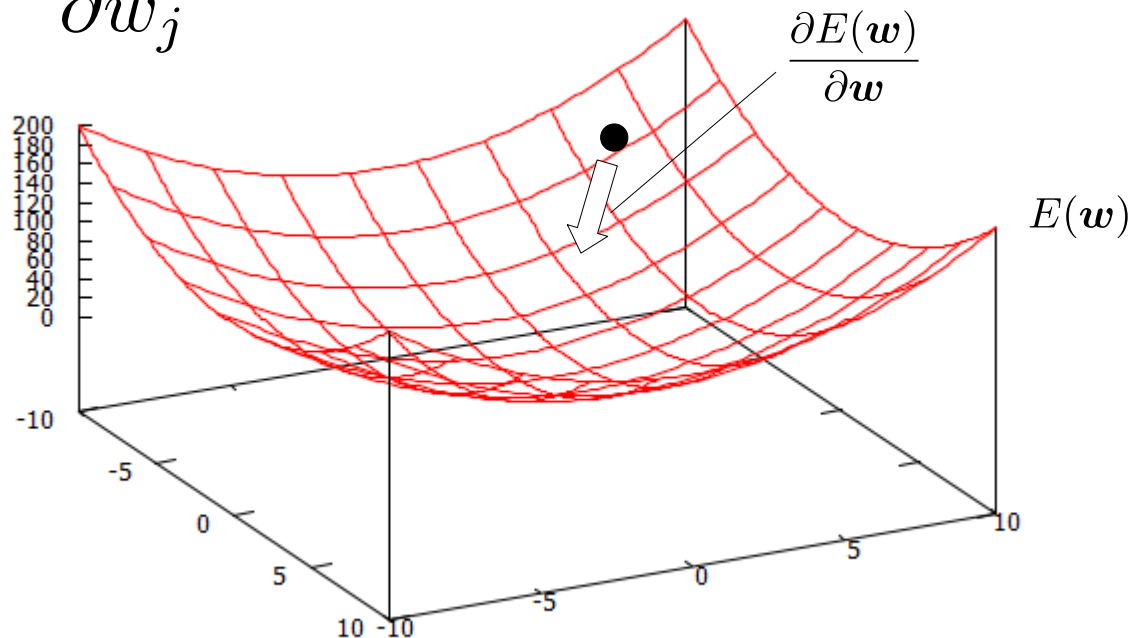
正解ラベル

$$= - \sum_{\boldsymbol{x}_i \in D} \{y_i \log o_i + (1 - y_i) \log(1 - o_i)\}$$

5.3.4 ロジスティック識別

- $E(w)$ を最急勾配法で最小化
 - 適当な初期値 w を選ぶ
 - w を $E(w)$ の勾配の逆方向に少しずつ修正

$$w_j \leftarrow w_j - \eta \frac{\partial E(w)}{\partial w_j}$$



5.3.4 ロジスティック識別

- 重み更新量の計算

$$\begin{aligned}\frac{\partial E(\boldsymbol{w})}{\partial w_j} &= \sum_{\boldsymbol{x}_i \in D} \frac{\partial E(\boldsymbol{w})}{\partial o_i} \cdot \frac{\partial o_i}{\partial w_j} \\ &= \sum_{\boldsymbol{x}_i \in D} \left(\frac{y_i}{o_i} - \frac{1 - y_i}{1 - o_i} \right) o_i (1 - o_i) x_{ij} \\ &= \sum_{\boldsymbol{x}_i \in D} (y_i - o_i) x_{ij}\end{aligned}$$

$E(\boldsymbol{w})$ を o_i で偏微分

シグモイド関数の微分

w_j の係数

- 重みの更新式

$$w_j \leftarrow w_j - \eta \sum_{\boldsymbol{x}_i \in D} (y_i - o_i) x_{ij}$$

5.3.5 確率的最急勾配法

- 最急勾配法の問題点
 - 全データに対して誤差を計算するので、データ数が多い場合、重み更新に時間がかかる
- 確率的最急勾配法
 - 個々のデータの計算結果に基づき重みを更新
 - データが来る毎に学習するオンライン学習が可能
- ミニバッチ法
 - 数十～数百程度のデータで誤差を計算し、修正方向を決める方法

学習のシミュレーションサイト

<https://playground.tensorflow.org/>

