

# 機械学習 第7回 サポートベクトルマシン

立命館大学 情報理工学部

福森 隆寛

Beyond Borders

# 講義スケジュール

(第1～4回、第14回) (第5～13回、第15回)

□ 担当教員：村上 陽平先生・福森 隆寛

|   |                |
|---|----------------|
| 1 | 機械学習とは、機械学習の分類 |
| 2 | 機械学習の基本的な手順    |
| 3 | 識別（１）          |
| 4 | 識別（２）          |
| 5 | 識別（３）          |
| 6 | 回帰             |
| 7 | サポートベクトルマシン    |
| 8 | ニューラルネットワーク    |

|    |           |
|----|-----------|
| 9  | 深層学習      |
| 10 | アンサンブル学習  |
| 11 | モデル推定     |
| 12 | パターンマイニング |
| 13 | 系列データの識別  |
| 14 | 強化学習      |
| 15 | 半教師あり学習   |

□ 担当教員：叶 昕辰先生（第16回の講義を担当）

# 講義スケジュール

## □ 担当教員 1 : 福森 (第1回～第15回)

|   |                    |
|---|--------------------|
| 1 | 機械学習とは、機械学習の分類     |
| 2 | 機械学習の基本的な手順        |
| 3 | 識別 ( 1 )           |
| 4 | 識別 ( 2 )           |
| 5 | 識別 ( 3 )           |
| 6 | 回帰                 |
| 7 | <b>サポートベクトルマシン</b> |
| 8 | ニューラルネットワーク        |

|    |           |
|----|-----------|
| 9  | 深層学習      |
| 10 | アンサンブル学習  |
| 11 | モデル推定     |
| 12 | パターンマイニング |
| 13 | 系列データの識別  |
| 14 | 半教師あり学習   |
| 15 | 強化学習      |

## □ 担当教員 2 : 叶昕辰先生 (第16回の講義を担当)

# 今回の講義内容

---

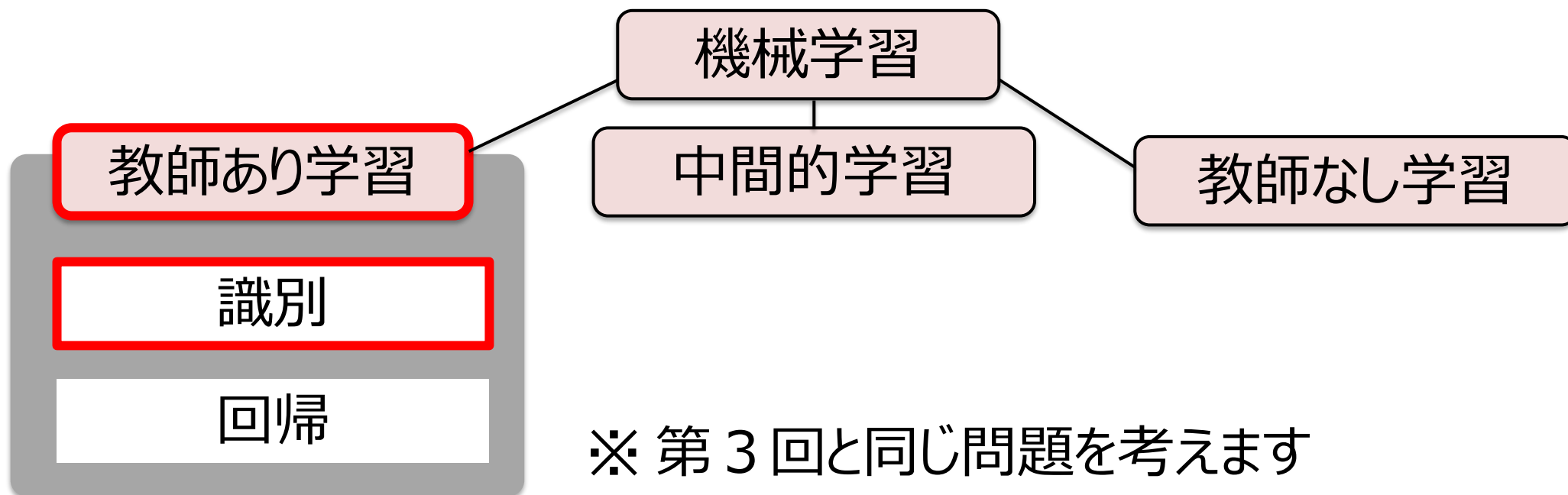
- 取り扱う問題の定義
- サポートベクトルマシン
  - ハードマージン
  - ソフトマージン
  - カーネル関数
- 演習問題

# 取り扱う問題の定義：教師あり・識別問題

- カテゴリデータ、または数値データからなる特徴ベクトルを入力して、それをクラス分けする識別器を作る

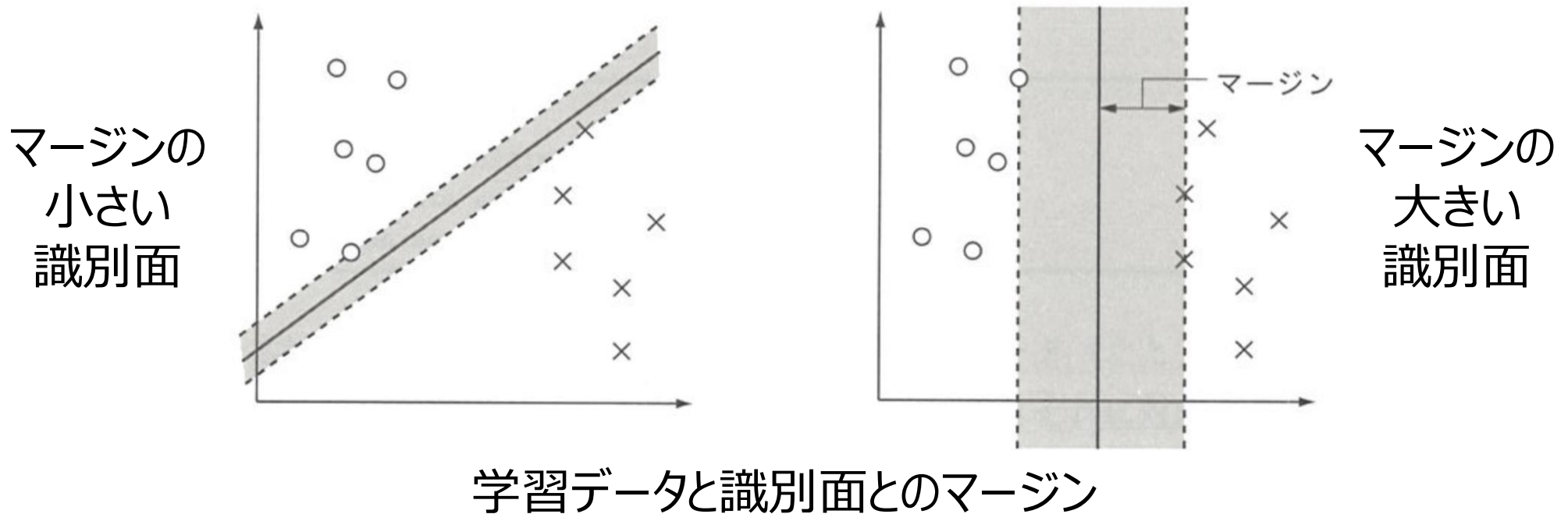
※ 教師あり学習の識別問題での学習データは、以下のペアで構成される

入力データの特徴ベクトル  $\leftarrow \{ \underline{x_i}, \underline{y_i} \}, \quad i = 1, 2, \dots, \underline{N} \longrightarrow$  学習データの総数  
(カテゴリデータ/数値データ) カテゴリ形式の正解情報  $\rightarrow$  「クラス」と呼ぶ



# サポートベクトルマシン

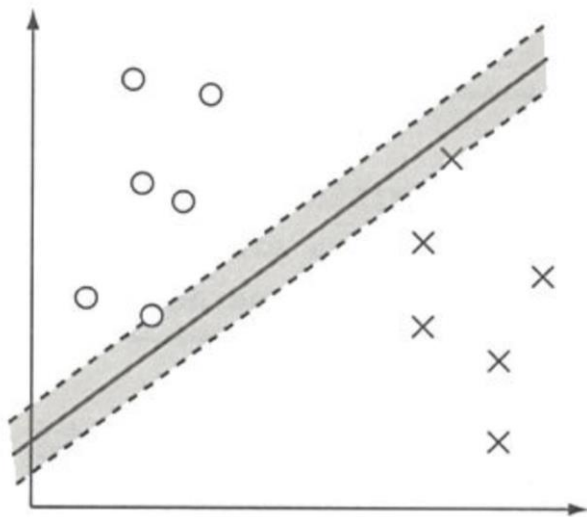
- 下図のような特徴空間上で超平面によって分離できる学習データを考える
  - 正解率100%で識別できる識別面が<sup>むすう</sup>無数に存在
  - **マージン**（識別面と最も近いデータとの距離）に基づいて識別面を決定（→汎化能力の高さを<sup>ていりょう</sup>定量的に表現）



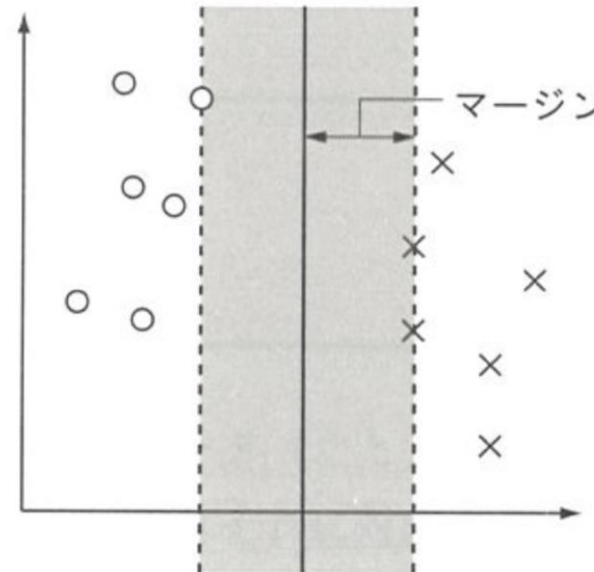
# サポートベクトルマシン

## □ サポートベクトルマシン (support vector machine: SVM)

- 学習データからのマージンが最大となる  
識別面（一般には識別超平面）を求める手法
- マージンが<sup>ひろ</sup>広いと、学習データと識別面の間に  
未知データが入る余地がある<sub>よ ち</sub>



マージンの小さい識別面



マージンの大きい識別面

学習データから  
少しずれた  
未知データが  
別クラスに  
識別されることが  
少なくなる

# サポートベクトルマシン：定式化

□ マージンが最大となる識別面を求める方法を考える

□ 前提条件

■ 学習データが線形分離可能な状況を仮定

■ 数値特徴に対して正解情報の付いたデータを使用

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

■ 2値分類問題に限定

• 正解情報 $y_i$ ：正例  $\rightarrow y_i = 1$ 、負例  $\rightarrow y_i = -1$



# サポートベクトルマシン：定式化

## □ 特徴空間上での識別面（超平面を仮定）

- 識別面： $\mathbf{w} \cdot \mathbf{x} + w_0 = 0$  ※  $\mathbf{w}$ 、 $\mathbf{x}$ ： $d$ 次元ベクトル
- 正例： $\mathbf{w} \cdot \mathbf{x} + w_0 > 0$
- 負例： $\mathbf{w} \cdot \mathbf{x} + w_0 < 0$

## □ 点と直線の距離の公式<sup>こうしき</sup>より

$i$ 番目のデータ $\mathbf{x}_i$ と、この識別面との距離  $\text{Dist}(\mathbf{x}_i)$

$$\text{Dist}(\mathbf{x}_i) = \frac{|\mathbf{w} \cdot \mathbf{x}_i + w_0|}{\|\mathbf{w}\|}$$

# サポートベクトルマシン：定式化

□ 識別面  $w \cdot x + w_0 = 0$  の両<sup>りょうへん</sup>辺を定数倍しても表す識別面は変わらない

□ 識別面に最も近いデータを識別面の式に代入したときその絶対値が1になるように係数 $w$ と $w_0$ を調整

$$\min_{i=1,\dots,N} |w \cdot x_i + w_0| = 1$$

この状況での学習パターンと識別面との最小距離は

$$\min_{i=1,\dots,N} \text{Dist}(x_i) = \min_{i=1,\dots,N} \frac{|w \cdot x_i + w_0|}{\|w\|} = \frac{1}{\|w\|}$$

この最小距離は「マージン」を表している

# サポートベクトルマシン：定式化

$$\min_{i=1,\dots,N} \text{Dist}(\mathbf{x}_i) = \min_{i=1,\dots,N} \frac{|\mathbf{w} \cdot \mathbf{x}_i + w_0|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

このマージンを最大にする識別面を求める問題は、 $\|\mathbf{w}\|$ を最小化する問題となる

□ これを $\|\mathbf{w}\|^2$ の最小化問題に置き換える

- 普通は  $\mathbf{w} = \mathbf{0}$  が最小解だが、これでは識別面にならない
- 識別面として全学習データを識別できるという条件を追加

$$y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \quad i = 1, \dots, N$$

- $y_i = 1$  または  $y_i = -1$  としていたので、  
正例・負例両方の制約を一つの式で表現

# サポートベクトルマシン：定式化

- マージンを最大にする識別面を求める問題の最終的な定式化

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{条件： } y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \quad i = 1, \dots, N$$

- <sup>のち</sup>後ほど、微分を利用して極値を求めて最小解を導くので  
上式のように乗数  $\frac{1}{2}$  を付けておく  
じょうすう

# サポートベクトルマシン：識別面の計算

□ 前スライドで定式化した問題を  
ラグランジュの未定乗数法を用いて解決

## □ ラグランジュ未定乗数法

■  $g(x) = 0$  の条件下で  $f(x)$  の最小値（または最大値）を  
計算する手法

■ ラグランジュ関数  $L(x, \lambda) = f(x) - \lambda g(x)$  を導入して、  
この関数の極値を求める問題に置き換える

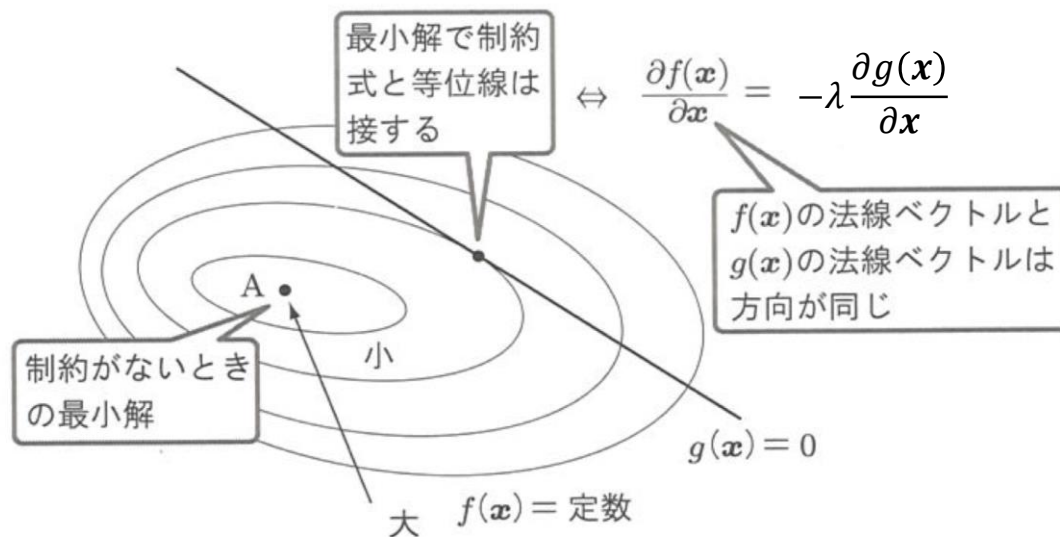
•  $\lambda$ ：ラグランジュ乗数

■ この  $x$  に関する偏微分を0とすると、下式が得られる

$$\frac{\partial f(x)}{\partial x} = -\lambda \frac{\partial g(x)}{\partial x}$$

# サポートベクトルマシン：識別面の計算

- 制約がない状況での最小解は、下図の点A
- 直線  $g(x) = 0$  の上で、 $f(x)$  が最低となる点  
→  $f(x)$  の等位線と直線  $g(x) = 0$  が接する点



この等位線と直線が接すること  
||  
それらの法線ベクトルが一致すること

$$\frac{\partial f(x)}{\partial x} = -\lambda \frac{\partial g(x)}{\partial x}$$

# サポートベクトルマシン：識別面の計算

(つづき)

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{条件: } y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \quad i = 1, \dots, N$$

の制約付きの最小化問題は、ラグランジュ乗数 $\alpha_i$ を導入して、関数 $L$ の最小値を求める問題に置き換える

$$L(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i \{y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) - 1\}$$

最小値では、 $L$ の勾配が0になるので、以下の式が成立<sup>せ い り つ</sup>

$$\frac{\partial L}{\partial w_0} = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0 \quad \frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

# サポートベクトルマシン：識別面の計算

(つづき)

$$L(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i \{y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) - 1\}$$

$$\sum_{i=1}^N \alpha_i y_i = 0, \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

これらの式を整理すると

$$L(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j = L_{\text{dual}}(\boldsymbol{\alpha})$$

「 $L(\mathbf{w}, w_0, \boldsymbol{\alpha})$ を最小にする問題」は「上に凸<sup>とつ</sup>の関数である  
 $L_{\text{dual}}(\boldsymbol{\alpha})$ を最大にする $\boldsymbol{\alpha}$ を計算する問題」と等価<sup>とうか</sup>



# サポートベクトルマシン：識別面の計算

- $L_{\text{dual}}(\alpha)$ が最大となる $\alpha$ を計算すると
  - $\alpha_i \neq 0$ ：サポートベクトルに対応するもののみ
  - $\alpha_i = 0$ ：サポートベクトル以外
  - サポートベクトル：識別面の計算に寄与する学習データ

- マージンを最大にする識別面の重み

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i, \quad w_0 = -\frac{1}{2} (\mathbf{w} \cdot \mathbf{x}_+ + \mathbf{w} \cdot \mathbf{x}_-)$$

- $\alpha_i$ ：  $L_{\text{dual}}(\alpha)$ が最大となる $\alpha$ の $i$ 番目の要素ばんめ
- $\mathbf{x}_+$ 、 $\mathbf{x}_-$ ：正例、負例に属するサポートベクトル

# サポートベクトルマシン：ソフトマージン

じっさい  
□ 実際は、線形分離不可能な学習データが多い

- ある程度の誤識別のデータを許容<sup>きょよう</sup>して、  
それらが識別面からあまり離れていない識別面を選択
- ソフトマージンの条件を導入

□ ハードマージン：全データを正しく識別できる条件

$$y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \quad i = 1, \dots, N$$

□ ソフトマージン：ハードマージンを弱<sup>よわ</sup>める変数 $\xi_i$ を導入

- $\xi_i (\geq 0)$ ： $i$ 番目のデータが制約を満たしていない程度で、  
小さい方がよい

$$y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1 - \xi_i, \quad i = 1, \dots, N$$

# サポートベクトルマシン：ソフトマージン

- SVMのマージン最大化問題に $\xi_i$ を導入

$$\min \left( \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \right)$$

$$y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1 - \xi_i, \quad i = 1, \dots, N, \quad \xi_i \geq 0$$

- 上式の $C$ は制約を満たさないデータを、どの程度の重みで最適化に組み込むかを決定する定数

- $C$ が大きい

- ハードマージンの問題設定に近づき、複雑な識別面になる

- $C$ が小さい

- 誤りをほぼ無視する振る舞いで、比較的単純な識別面になる

# サポートベクトルマシン：ソフトマージン

- 前スライドをラグランジュの未定乗数法で解くと、ハードマージンでの識別面と同じ式が導出され、

$$L_{\text{dual}}(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

が最大となる $\boldsymbol{\alpha}$ を計算する

ただし、ソフトマージンでは  $0 \leq \alpha_i \leq C$  の制約がつく

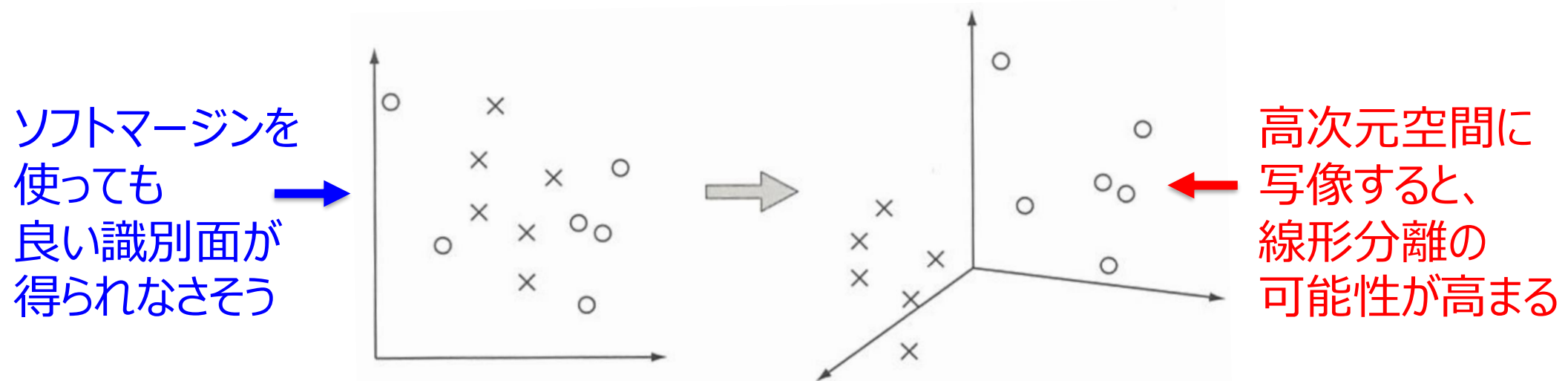
# 演習問題7-1（10分間）

---

- SVMは2クラスのデータに対する分類器である。  
ここで、3クラス以上のデータに対する多クラス  
識別問題にSVMを適用するには、  
どうすれば良いか考えよ

# サポートベクトルマシン：カーネル関数

- 信頼できる識別面が得られそうにない場合を考える
- 特徴空間の次元数 $d$ が大きい場合
  - データが高次元空間上に疎らに分布することになるので、線形識別面が存在する可能性が高くなる



低次元の特徴ベクトルを高次元空間に写像し、  
高次元空間上でSVMを使って識別超平面を求める方法を考える

# サポートベクトルマシン：カーネル関数

- $d$ 次元の特徴空間に対して  
元の空間におけるデータ間の距離関係を  
保存しながら高次元に非線形に写像する
  - 高次元空間上での線形識別器の性能は  
元の空間での複雑な非線形識別器の性能に相当する
  - ただし、識別に無関係な特徴を持ち込むと、データが  
無意味な方向に疎らに分布し、元の分布の性質が  
壊されやすくなる
- 元の空間におけるデータ間の距離関係を  
保存するような非線形写像が見つかるか？

こわ

# サポートベクトルマシン：カーネル関数

□ 元の特徴空間上の2点 $x, x'$ の距離に基づいて定義される類似度関数 $K(x, x')$ を考える

■ 2点 $x, x'$ が近いほど、 $K(x, x')$ は大きな値になる

□  $K(x, x')$ が、半正定値性などの条件を満たすと

$$K(x, x') = \phi(x)^T \phi(x')$$

のように、2点 $x, x'$ から求まる高次元ベクトルの内積ないせきによって、類似度関数の値を計算できる

■ この関数を**カーネル関数**と呼ぶ



# サポートベクトルマシン：カーネル関数

## □ カーネル関数の例

### ■ 多項式カーネル

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^p \quad \text{※ } p : \text{自然数}$$

### ■ ガウシアンカーネル関数

$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

- $\gamma$  : カーネル関数の広がりを表すパラメータ  $\sigma^2$  の<sup>ぎやくすう</sup>逆数  $(\frac{1}{\sigma^2})$
- $\gamma$  が大きい場合は、複雑な識別面が形成される
  - カーネル関数が大きな値をとる範囲が<sup>せま</sup>狭くなるので、  
識別面の形成に<sup>けいせい</sup>関与するデータが<sup>かんよ</sup>近傍の<sup>かぎ</sup>ものに限られる

# サポートベクトルマシン：カーネル関数

- 写像後の空間での識別関数  $g(\phi(x))$

$$g(\phi(x)) = \mathbf{w}^T \phi(x) + w_0$$

- ここでSVMを適用すると、 $\mathbf{w}$ は

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \phi(\mathbf{x}_i)$$

となるので、識別関数  $g(\phi(x))$  は

$$\begin{aligned} g(\phi(x)) &= \sum_{i=1}^N \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(x) + w_0 \\ &= \sum_{i=1}^N \alpha_i y_i K(x, \mathbf{x}_i) + w_0 \end{aligned}$$

# サポートベクトルマシン：カーネル関数

- これまでと同様に、学習の問題も

$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

が最大となる $\alpha$ を計算すれば良い

- カーネル関数 $K$ さえ定まれば、  
識別面を得られるのが重要なポイント

## □ カーネルトリック

- 複雑な非線形写像を求める操作を避ける方法
- これがSVMが色々な応用に使われてきた理由

# サポートベクトルマシン：カーネル関数

## □ カーネル法

- 入力データを高次元空間に写像しながら、計算上は明示的に高次元空間を考えずに識別面を構成できる

## □ 非線形写像で線形分離可能な高次元空間にデータを飛ばして、マージン最大化基準で信頼できる識別面を求めるSVMは非常に強力

- 文書分類やバイオインフォマティクスなどの様々な分野（特に高次元識別問題）で利用

# 演習問題7-2（10分間）

- カーネル関数の1つである多項式カーネルは、以下の式で与えられる

$$K(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{\phi}(\boldsymbol{x})^T \boldsymbol{\phi}(\boldsymbol{x}') = (\boldsymbol{x}^T \boldsymbol{x}' + 1)^p$$

- 2次元特徴ベクトルを  $\boldsymbol{x} = (x_1, x_2)$ 、 $p = 2$ として  $\boldsymbol{\phi}(\boldsymbol{x})$ を求めよ

# 補足資料

サポートベクトルマシンのマージンを最大とする識別面の計算にて

$$L(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i \{y_i \cdot (\mathbf{w} \cdot \mathbf{x}_i + w_0) - 1\}$$

$$\sum_{i=1}^N \alpha_i y_i = 0, \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

から

$$L(\mathbf{w}, w_0, \boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j = L_{\text{dual}}(\boldsymbol{\alpha})$$

を求める

# 補足資料

$$\begin{aligned} L(\mathbf{w}, w_0, \alpha) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i \{y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) - 1\} \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i y_i (\mathbf{w} \cdot \mathbf{x}_i + w_0) + \sum_{i=1}^N \alpha_i \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \mathbf{w} \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i - w_0 \sum_{i=1}^N \alpha_i y_i + \sum_{i=1}^N \alpha_i \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \mathbf{w} \cdot \mathbf{w} - w_0 \cdot 0 + \sum_{i=1}^N \alpha_i \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i = -\frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i \end{aligned}$$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$= -\frac{1}{2} \left\| \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right\|^2 + \sum_{i=1}^N \alpha_i$$

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

# 補足資料

(つづき)

$$\begin{aligned} L(\mathbf{w}, w_0, \boldsymbol{\alpha}) &= -\frac{1}{2} \left\| \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right\|^2 + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} \left( \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)^T \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right) + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} (\alpha_1 y_1 \mathbf{x}_1^T + \cdots + \alpha_N y_N \mathbf{x}_N^T) \left( \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right) + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} \left( \alpha_1 y_1 \mathbf{x}_1^T \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j + \cdots + \alpha_N y_N \mathbf{x}_N^T \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right) + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} \left( \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i^T \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right) + \sum_{i=1}^N \alpha_i = -\frac{1}{2} \left( \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right) + \sum_{i=1}^N \alpha_i \end{aligned}$$