

教師なし学習

佐藤 一誠

sato@k.u-tokyo.ac.jp

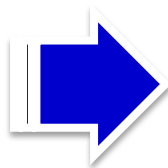
<http://www.ms.k.u-tokyo.ac.jp>

教師なし学習とは

2

学習データ:

$$\{x_i\}_{i=1}^n$$



$$x_i = f_{\theta}(z_i)$$

データに共通する部分

データに固有な部分

教師付き学習は

$$\{y_i, x_i\}_{i=1}^n$$



$$y_i = f_{\theta}(x_i)$$

講義の流れ



1. 次元削減
2. クラスタリング
3. 非線形化
4. 生成モデル

データに共通する部分

$$x_i = f_{\theta}(z_i)$$
The equation $x_i = f_{\theta}(z_i)$ is shown. The parameter θ is enclosed in a light green circle, and the latent variable z_i is enclosed in a light red circle. A green line points from the text 'データに共通する部分' to the green circle, and a red line points from the text 'データに固有な部分' to the red circle.

データに固有な部分



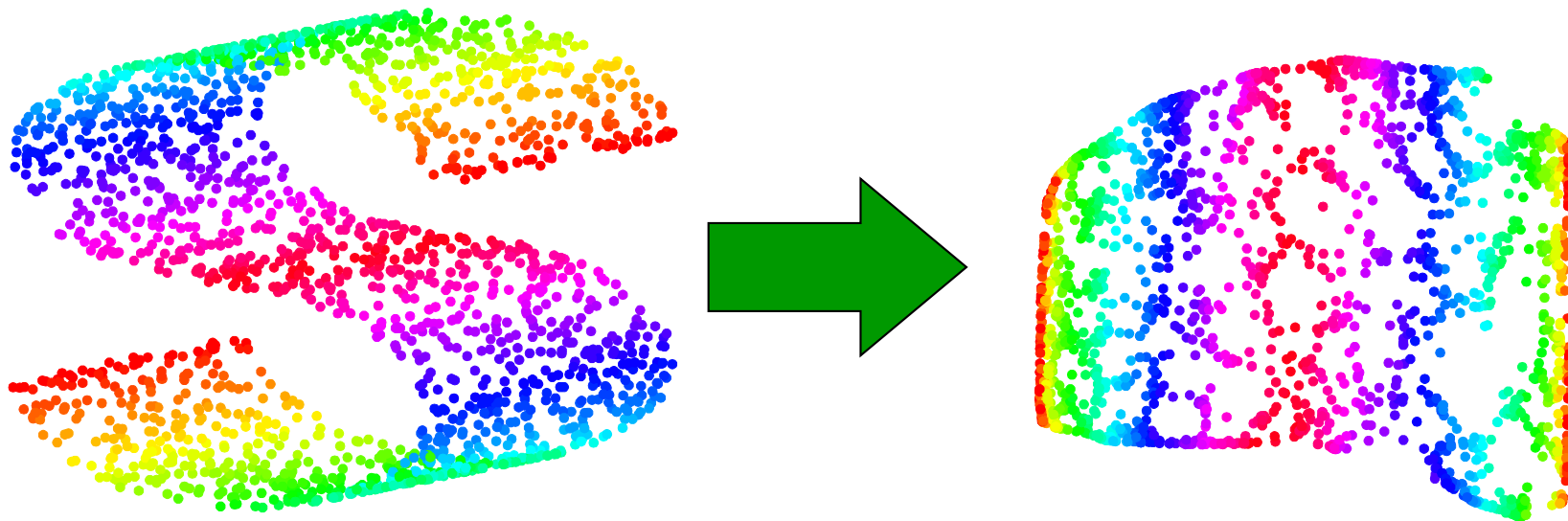
[z が連続の場合]
データの次元よりも
低い次元での表現

次元削減

- 本質的な情報を保持したまま次元を減らしたい！

もとの高次元データ

次元削減後のデータ(2d)



- 1～3次元に減らせば, データを可視化できる.
- 次元削減の基本的な仮定:
 - 手持ちの高次元データはある意味で冗長である

線形次元削減

■ (高次元) 標本:

$$\{\mathbf{x}_i\}_{i=1}^n, \quad \mathbf{x}_i \in \mathbb{R}^d, \quad d \gg 1$$

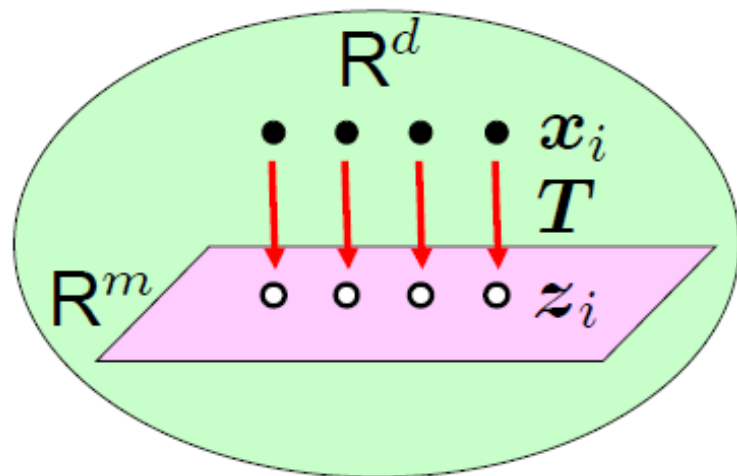
■ 埋め込み行列:

$$\mathbf{T} \in \mathbb{R}^{m \times d}, \quad 1 \leq m \ll d$$

■ 埋め込まれた標本

$$\{\mathbf{z}_i\}_{i=1}^n, \quad \mathbf{z}_i = \mathbf{T}\mathbf{x}_i \in \mathbb{R}^m$$

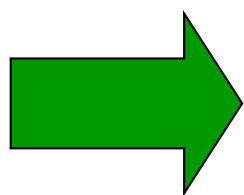
$$m \left\{ \begin{array}{c} \boxed{\mathbf{z}_i} \end{array} \right\} = \boxed{\mathbf{T}} \left\{ \begin{array}{c} \boxed{\mathbf{x}_i} \end{array} \right\} d$$



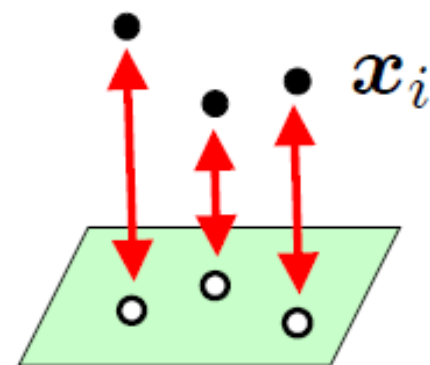
主成分分析 (PCA: Principal Component Analysis)

- **考え方**: データを表現するうえで最も重要な次元(部分空間)を取り出す

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 0.1 \end{pmatrix}, \begin{pmatrix} 3 \\ -0.1 \end{pmatrix}$$



低次元の部分空間に
正射影したときに
できるだけデータが
変化しないようにする

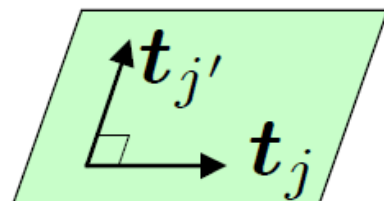


正射影

7

- $\{\mathbf{t}_j \mid \mathbf{t}_j \in \mathbb{R}^d\}_{j=1}^m$: m 次元部分空間の
正規直交基底

$$\mathbf{t}_j^\top \mathbf{t}_{j'} = \begin{cases} 1 & (j = j') \\ 0 & (j \neq j') \end{cases}$$

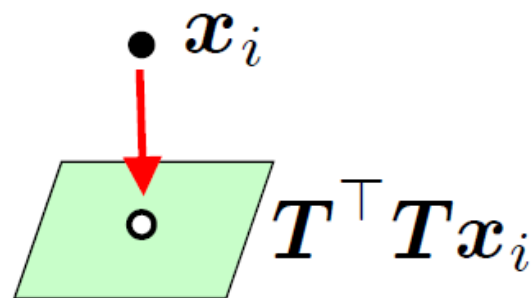


- 行列で表現すると, $\mathbf{T}\mathbf{T}^\top = \mathbf{I}_m$

$$\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_m)^\top \in \mathbb{R}^{m \times d}$$

- 標本 \mathbf{x}_i の正射影は

$$\sum_{j=1}^m (\mathbf{t}_j^\top \mathbf{x}_i) \mathbf{t}_j \quad \left(= \mathbf{T}^\top \mathbf{T} \mathbf{x}_i \right)$$



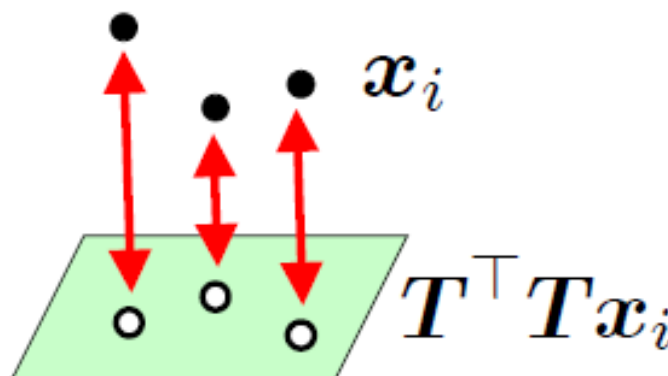
主成分分析の規準

- 射影誤差の和を最小にする:

$$\mathbf{T}_{\text{PCA}} = \underset{\mathbf{T} \in \mathbb{R}^{m \times d}}{\operatorname{argmin}} \left[\sum_{i=1}^n \|\mathbf{T}^\top \mathbf{T} \mathbf{x}_i - \mathbf{x}_i\|^2 \right]$$

$$\text{subject to } \mathbf{T}\mathbf{T}^\top = \mathbf{I}_m$$

(正規直交性の拘束条件をつける)



■ 次式を証明せよ

$$\sum_{i=1}^n \|T^{\top} T x_i - x_i\|^2 = -\text{tr} \left(T C T^{\top} \right) + \text{tr} (C)$$

$\text{tr}(\cdot)$: 行列のトレース

$$C = \sum_{i=1}^n x_i x_i^{\top}$$

● 仮定: $\frac{1}{n} \sum_{i=1}^n x_i = 0$

標本の散布行列
(正規化されていない
共分散行列)

■ ヒント: $T^{\top} T T^{\top} T = T^{\top} T$

$$\|y\|^2 = y^{\top} y$$

解答例

10

$$\sum_{i=1}^n \|T^\top T x_i - x_i\|^2$$

$$= \sum_{i=1}^n x_i^\top T^\top T T^\top T x_i - 2 \sum_{i=1}^n x_i^\top T^\top T x_i + \sum_{i=1}^n x_i^\top x_i$$

$T^\top T T^\top T = T^\top T$: 2回射影しても変わらない

$$= - \sum_{i=1}^n x_i^\top T^\top T x_i + \sum_{i=1}^n x_i^\top x_i$$

$$a^\top b = \text{tr}(ba^\top)$$

$$= - \sum_{i=1}^n \text{tr}(T x_i x_i^\top T^\top) + \sum_{i=1}^n \text{tr}(x_i x_i^\top)$$

$$= -\text{tr}(T C T^\top) + \text{tr}(C)$$

$$C = \sum_{i=1}^n x_i x_i^\top$$

固有値問題

- 対称行列 C は, 必ず

$$C = \sum_{j=1}^d \lambda_j \xi_j \xi_j^\top$$

という形に分解することができる.

- λ_j, ξ_j は固有値, 固有ベクトルとよばれ, 次の固有方程式の解として与えられる

$$C\xi = \lambda\xi$$

- 固有ベクトル ξ_1, \dots, ξ_d は互いに直交する:

$$\xi_j^\top \xi_{j'} = 0 \quad \text{for } j \neq j'$$

- Octaveではeig関数で計算できる

固有値問題

$$C = \sum_{j=1}^d \lambda_j \boldsymbol{\xi}_j \boldsymbol{\xi}_j^\top$$

- 固有値分解を用いれば, 逆行列は次式で,

$$C^{-1} = \sum_{j=1}^d \lambda_j^{-1} \boldsymbol{\xi}_j \boldsymbol{\xi}_j^\top$$

$$\forall \lambda_j \neq 0$$

行列の平方根は次式で求められる

$$C^{1/2} = \sum_{j=1}^d \lambda_j^{1/2} \boldsymbol{\xi}_j \boldsymbol{\xi}_j^\top$$

$$\forall \lambda_j > 0$$

主成分分析の解

$$\mathbf{T}_{\text{PCA}} = \operatorname{argmax}_{\mathbf{T} \in \mathbb{R}^{m \times d}} \operatorname{tr} \left(\mathbf{T} \mathbf{C} \mathbf{T}^{\top} \right)$$

$$\text{subject to } \mathbf{T} \mathbf{T}^{\top} = \mathbf{I}_m$$

■ 主成分分析の解は次式で与えられる:

$$\mathbf{T}_{\text{PCA}} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_m)^{\top}$$

- $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_m$: は固有値問題 $\mathbf{C}\boldsymbol{\xi} = \lambda\boldsymbol{\xi}$ の固有値 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ に対応する正規化 $\|\boldsymbol{\xi}_j\| = 1$ された固有ベクトル

主成分分析の解の求め方

14

1. 固有値問題を解く:

標本の散布行列

$$C\xi = \lambda\xi$$

$$C = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

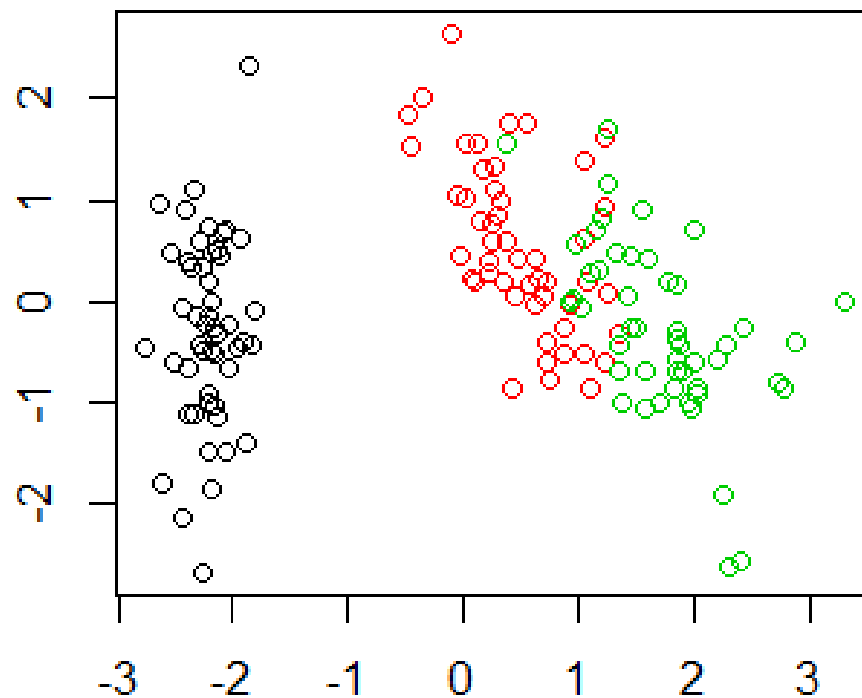
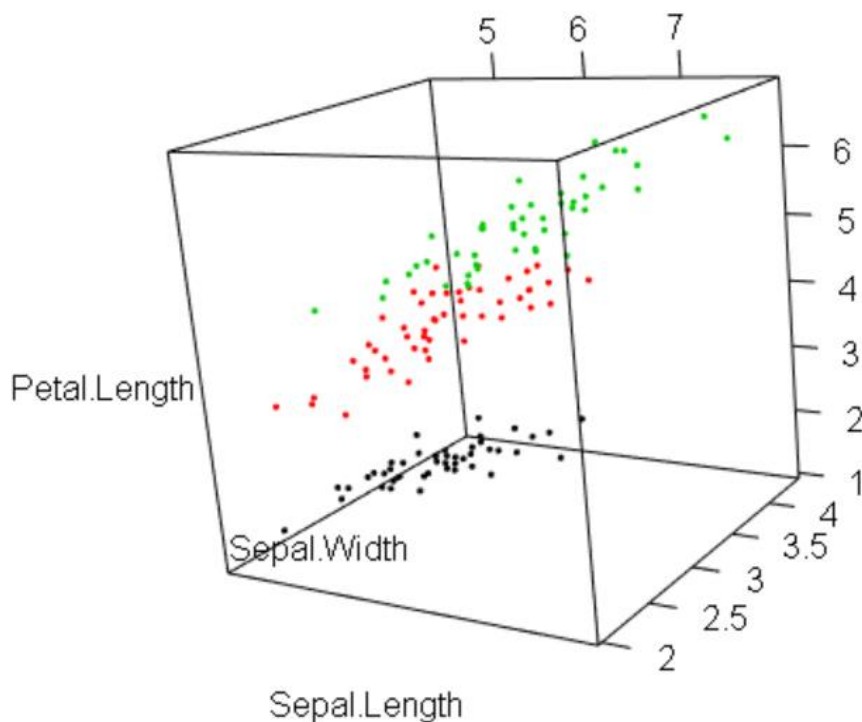
- 固有値を降順にソート: $\lambda_1 \geq \dots \geq \lambda_d$
- 固有ベクトルを正規化: $\|\xi_j\| = 1$

2. 上位 m 個の固有ベクトルを並べる:

$$T_{\text{PCA}} = (\xi_1, \dots, \xi_m)^\top$$

実行例(IRISデータ): 3D \Rightarrow 2D

15



- 主成分分析によって、データの大局的な分布をよく表す部分空間が得られている
- しかし、クラスタ構造のようなデータの局所的な構造は失われることがある

講義の流れ



1. 次元削減
2. クラスタリング
3. 非線形化
4. 生成モデル

データに共通する部分

$$x_i = f_{\theta}(z_i)$$
The equation $x_i = f_{\theta}(z_i)$ is shown. The parameter θ is enclosed in a light green circle, and the latent variable z_i is enclosed in a light red circle. A green line connects the text 'データに共通する部分' to the green circle, and a red line connects the text 'データに固有な部分' to the red circle.

データに固有な部分

[zが離散の場合]

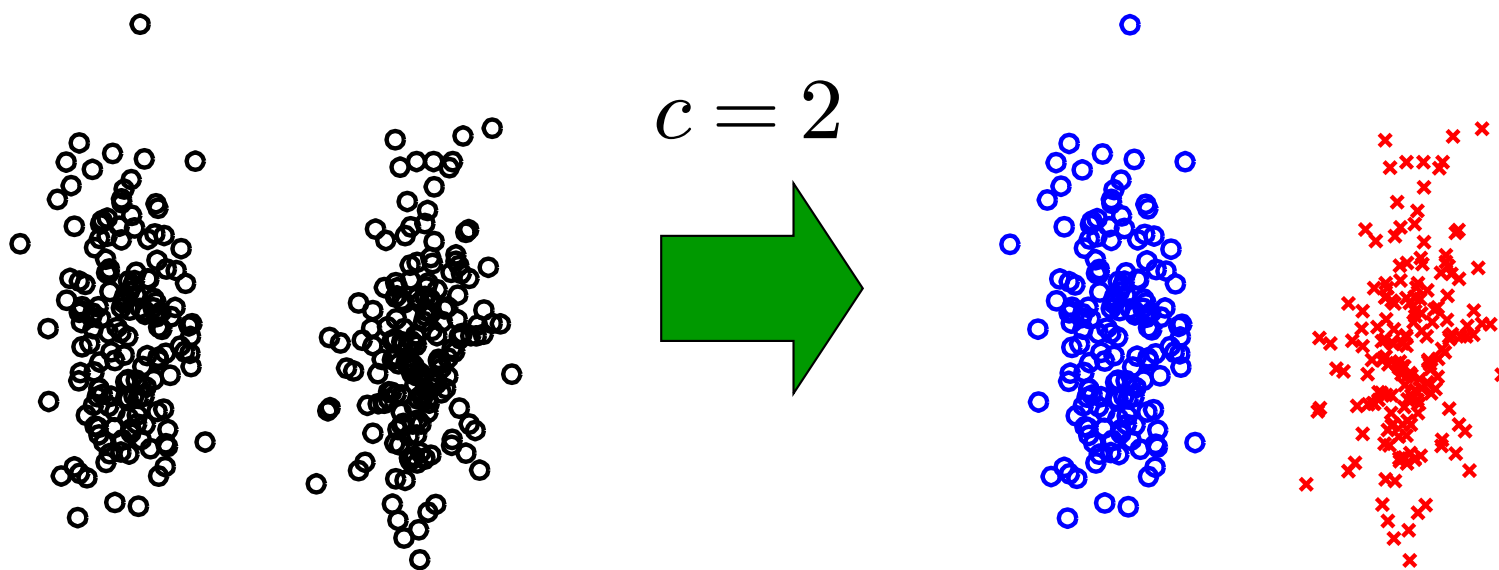
データのクラスタ
構造を表現



データのクラスタリング

17

- **目標**: ラベル無しの標本 $\{x_i\}_{i=1}^n$ を c 個のグループに分ける
 - 似た性質を持つデータは同じグループに
 - 異なる性質を持つデータは違うグループに
- c はあらかじめ固定しておく.



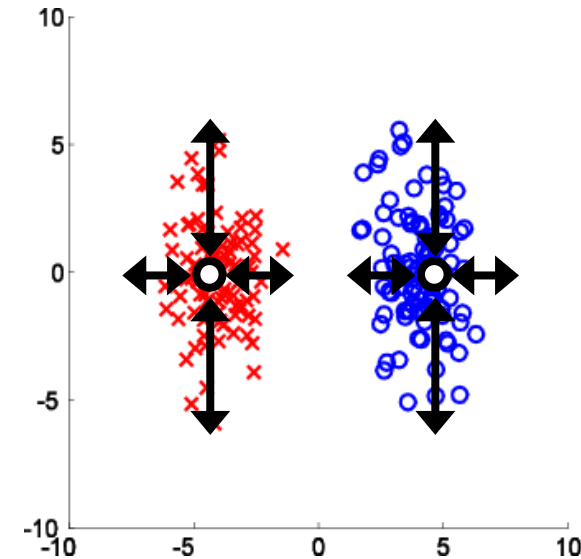
クラスタ内散布和の最小化

18

■ クラスタ内散布 :

$$\sum_{i:y_i=y} \|x_i - \mu_y\|^2$$

$$\mu_y = \frac{1}{n_y} \sum_{i:y_i=y} x_i$$



■ 考え方 : クラスタ内散布の和を最小にするように標本をクラスタに割り当てる

$$\min_{y_1, \dots, y_n \in \{1, \dots, c\}} \sum_{y=1}^c \sum_{i:y_i=y} \|x_i - \mu_y\|^2$$

■ しかしこれはNP困難 (計算量が $O(c^n)$) なので実時間では解けない.

1. クラスタ中心 $\{\mu_y\}_{y=1}^c$ を(ランダムに)初期化する
2. 収束するまで以下を繰り返す:

A) 各標本のクラスタの割り当てを次式で更新する:

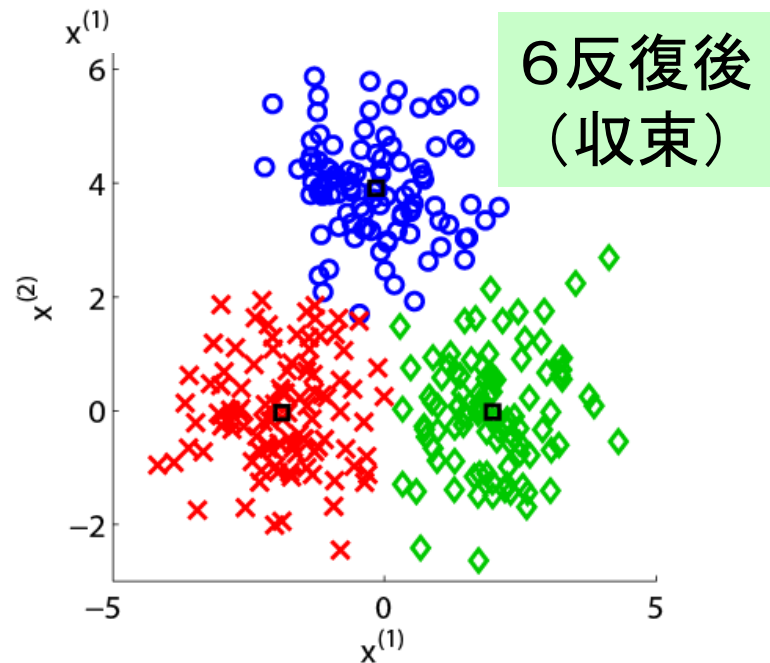
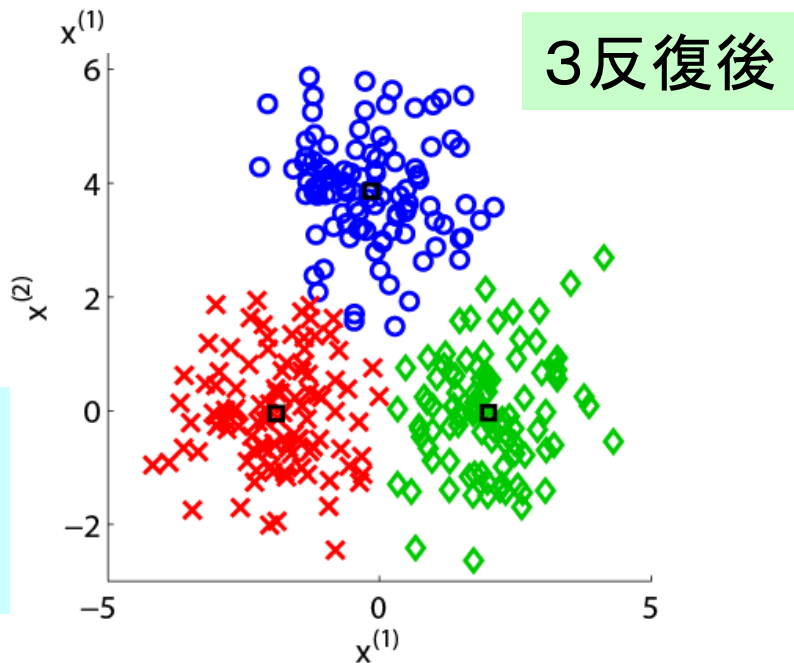
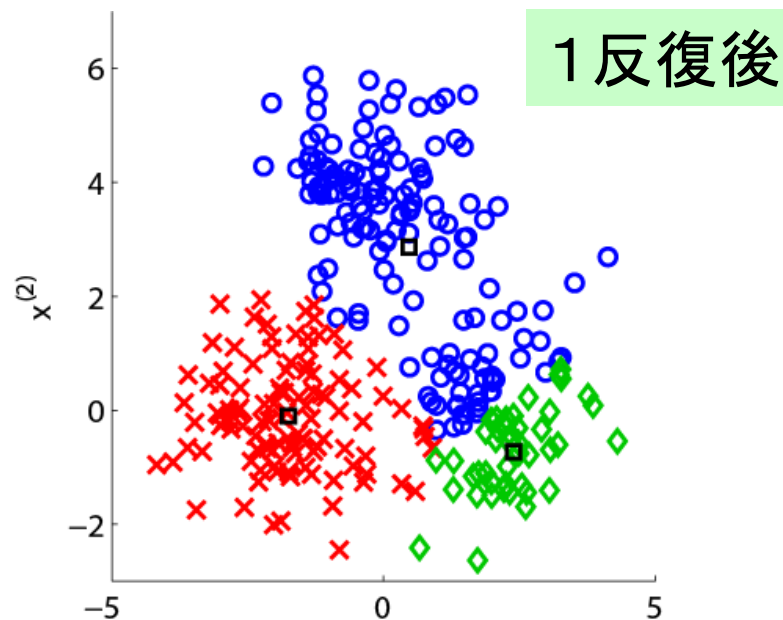
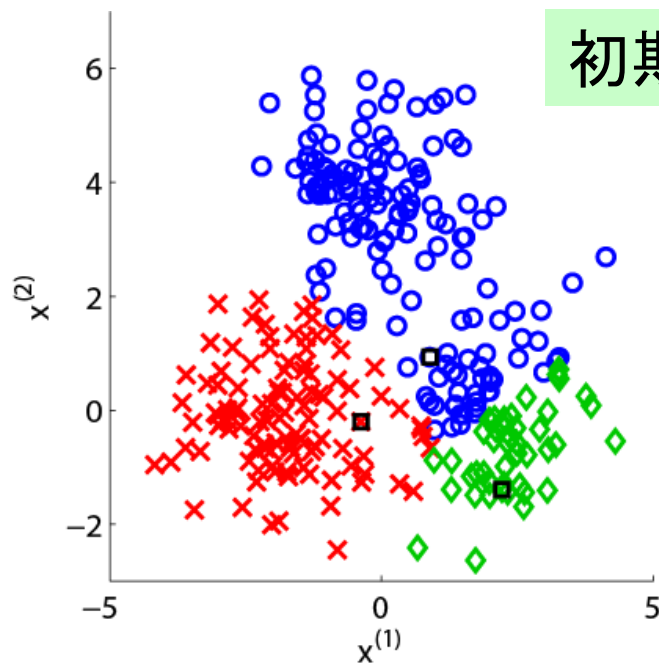
$$y_i \longleftarrow \operatorname{argmin}_{y \in \{1, \dots, c\}} \|\mathbf{x}_i - \mu_y\|^2, \quad i = 1, \dots, n$$

B) 各クラスタの中心を更新する:

$$\mu_y \longleftarrow \frac{1}{n_y} \sum_{i: y_i = y} \mathbf{x}_i, \quad y = 1, \dots, c$$

- これは, クラスタ内散布和の最小化問題の**局所最適解**を求めるアルゴリズムになっている.

実行例

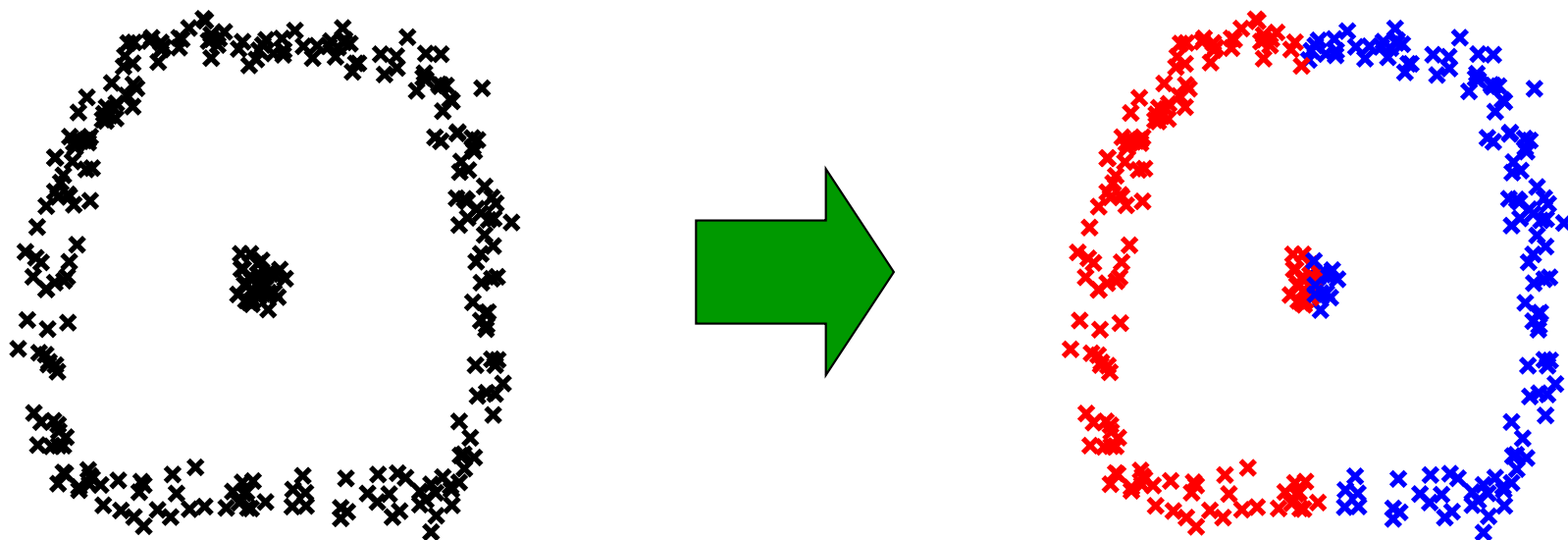


実装は
宿題

クラスタリング:まとめ

21

- クラスタ内散布の和を最小にする
- そのままではNP困難なので近似解を求める
 - クラスタリング結果が初期値に依存する.
- クラスタの形が凸でないとき, うまくいかない.
- クラスタ数をあらかじめ決める必要がある.



講義の流れ

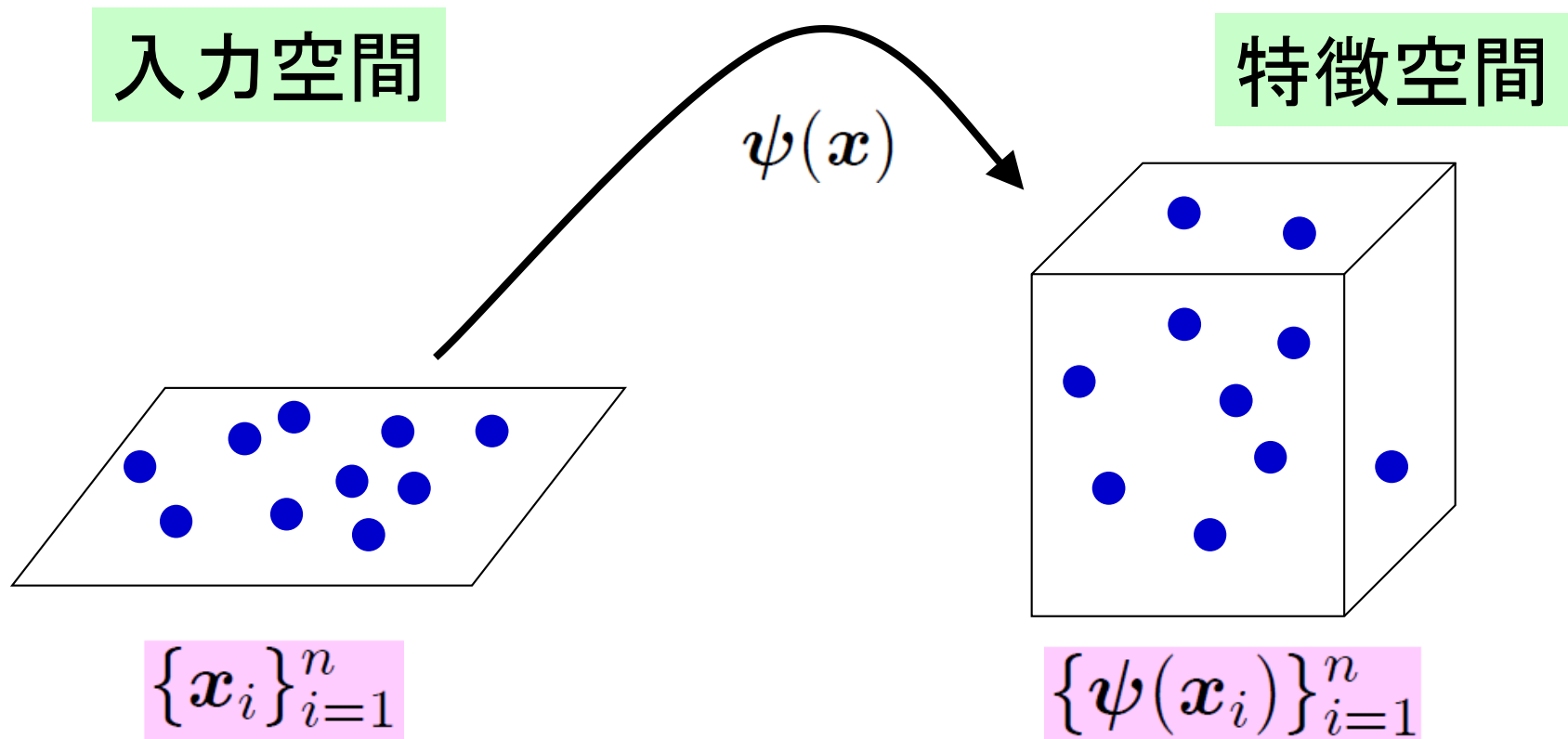


1. 次元削減
2. クラスタリング
3. 非線形化
4. 生成モデル

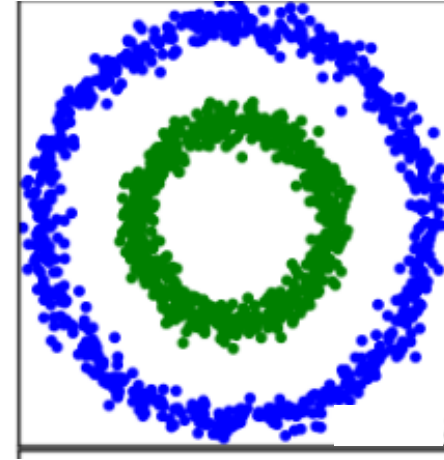
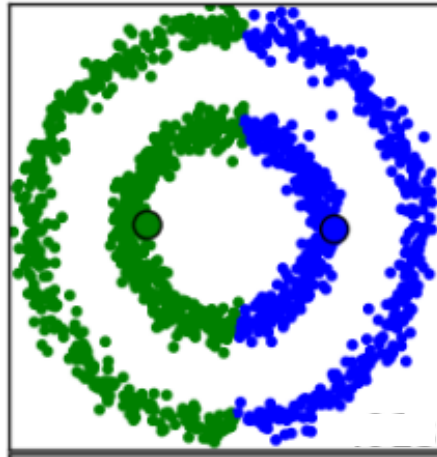
非線形への拡張

23

- 非線形関数 $\psi(x)$ で標本を特徴空間へ写像し、特徴空間内でアルゴリズムを実行



K-meansの非線形拡張



直感的には

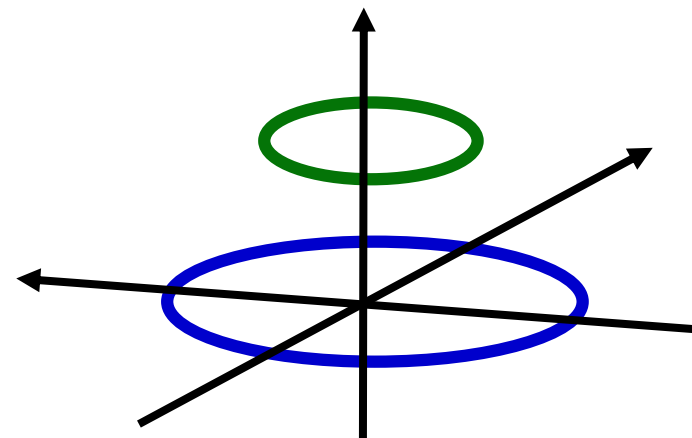
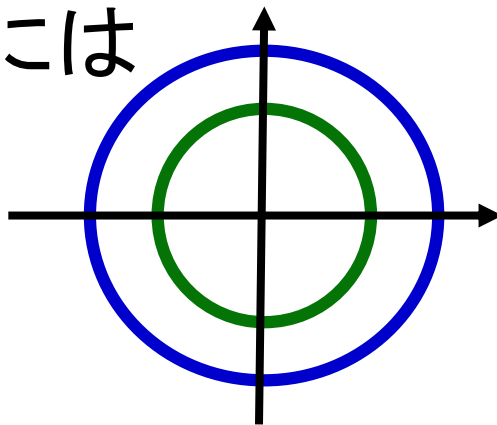


Image from http://scikit-learn.org/stable/auto_examples/cluster/plot_cluster_comparison.html

カーネルトリック

- 特徴空間内での内積を半正定値カーネル関数で計算:

$$\psi(x_i)^\top \psi(x_j) = K(x_i, x_j)$$

$$\forall x, x', K(x, x') \geq 0$$

例えばガウシアンカーネル

$$K(x, x') = \exp(-\|x - x'\|^2 / (2h^2))$$

- 主成分分析, k平均クラスタリングはカーネル非線形化できる
 - しかし, カーネルの選択が自明でない

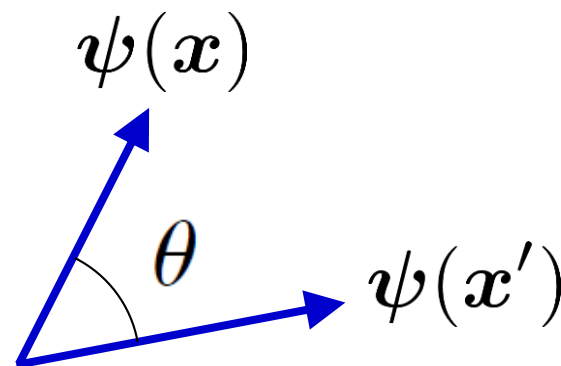
$$\psi(x)^\top \psi(x') = K(x, x')$$

■ 以下の量をカーネル関数を用いて計算せよ

● ノルム: $\|\psi(x)\|$

● 距離: $\|\psi(x) - \psi(x')\|$

● 角度: $\cos \theta$



$$\psi(x)^\top \psi(x') = \|\psi(x)\| \|\psi(x')\| \cos \theta$$

$$\blacksquare \|\psi(\mathbf{x})\| = \sqrt{\|\psi(\mathbf{x})\|^2} = \sqrt{\psi(\mathbf{x})^\top \psi(\mathbf{x})} = \sqrt{K(\mathbf{x}, \mathbf{x})}$$

$$\begin{aligned} \blacksquare \|\psi(\mathbf{x}) - \psi(\mathbf{x}')\| &= \sqrt{\|\psi(\mathbf{x}) - \psi(\mathbf{x}')\|^2} \\ &= \sqrt{\|\psi(\mathbf{x})\|^2 - 2\psi(\mathbf{x})^\top \psi(\mathbf{x}') + \|\psi(\mathbf{x}')\|^2} \\ &= \sqrt{K(\mathbf{x}, \mathbf{x}) - 2K(\mathbf{x}, \mathbf{x}') + K(\mathbf{x}', \mathbf{x}')} \end{aligned}$$

$$\blacksquare \cos \theta = \frac{\psi(\mathbf{x})^\top \psi(\mathbf{x}')}{\sqrt{\|\psi(\mathbf{x})\|^2 \|\psi(\mathbf{x}')\|^2}} = \frac{K(\mathbf{x}, \mathbf{x}')}{\sqrt{K(\mathbf{x}, \mathbf{x}) K(\mathbf{x}', \mathbf{x}')}}}$$

$$\psi(x)^\top \psi(x') = K(x, x')$$

- カーネル関数を用いてk-平均法の以下の式を表現せよ

$$y_i \longleftarrow \operatorname{argmin}_{y \in \{1, \dots, c\}} \|\mathbf{x}_i - \boldsymbol{\mu}_y\|^2, \quad i = 1, \dots, n$$

ヒント:

$$k = \arg \min_k \left\| \psi(x_i) - \frac{1}{n} \sum_{j: z_j = k} \psi(x_j) \right\|^2$$

$$\begin{aligned} & \left\| \psi(x_i) - \frac{1}{n} \sum_{j: z_j = k} \psi(x_j) \right\|^2 \\ &= \|\psi(x_i)\|^2 - \frac{2}{n} \psi(x_i)^T \sum_{j: z_j = k} \psi(x_j) + \frac{1}{n^2} \sum_{j: z_j = k} \sum_{j': z_{j'} = k} \psi(x_j) \psi(x_{j'}) \\ &= K(x_i, x_i) - \frac{2}{n} \sum_{j: z_j = k} K(x_i, x_j) + \frac{1}{n^2} \sum_{j: z_j = k} \sum_{j': z_{j'} = k} K(x_j, x_{j'}) \end{aligned}$$

講義の流れ



1. 次元削減
2. クラスタリング
3. 非線形化
4. 生成モデル

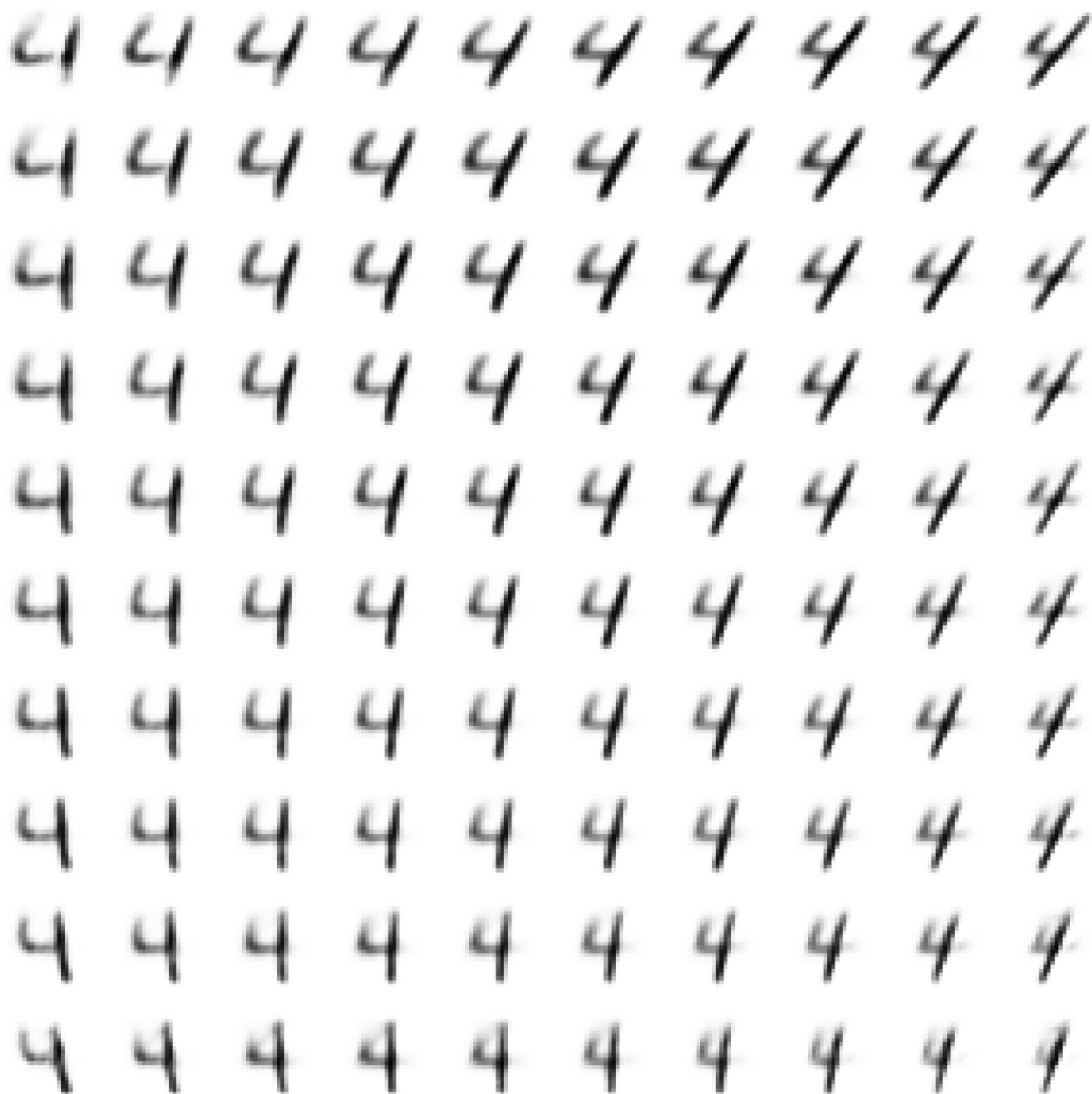
データに共通する部分

$$x_i = f_{\theta}(z_i)$$

データに固有な部分

zを入力すれば
xを生成できる





講義の流れ



1. 次元削減
2. クラスタリング
3. 非線形化
4. 生成モデル

データに共通する部分

$$x_i = f_{\theta}(z_i, y_i)$$

連続 離散

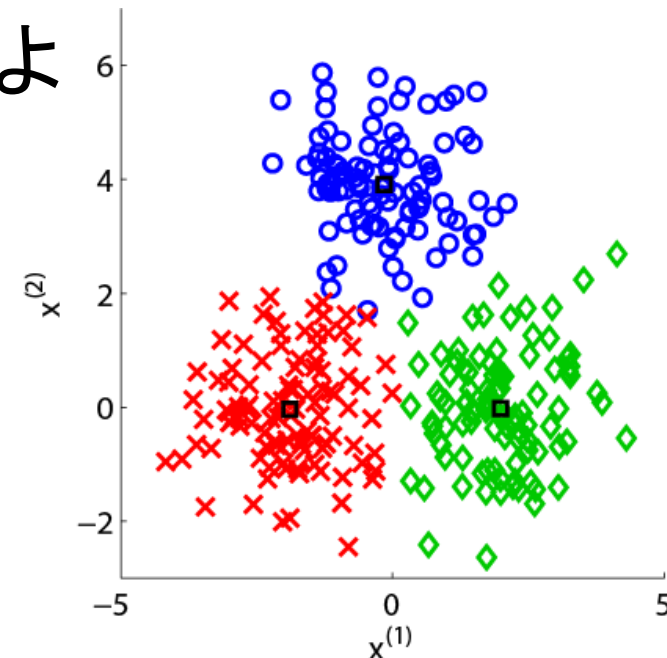
データに固有な部分連続と離散の両方を用いる
離散部分はカテゴリ情報があるデータ
に関しては予め入力しておく



文章の生成

「レンタ・カーは空のグラスを手にとり、蛇腹はすっかり暗くなっていた。それはまるで獲物を咀嚼しているようだった。彼は僕と同じようなものですね」と私は言った。「でもあなたはよく女の子に爪切りを買った。そしてその何かを振り払おうとしたが、今では誰にもできやしないのよ。私は長靴を棚の上を乗り越えるようにした。...

- 3次元以上のデータ(任意)に対してPCAを適用し2次元へ次元圧縮せよ(実装例を示す)
- 2次元に次元圧縮したデータを用いてk平均クラスタリングを適用せよ(実装例を示す)
- クラスタリング結果を可視化せよ



可視化例

- **教師なし学習**: 入力データだけから学習
 - **次元削減**: 高次元データに含まれる本質的な情報を保持したまま次元数を削減
 - **クラスタリング**: データをグループに分割
 - **生成モデル**: 訓練データから新たなデータを生成
- **カーネルトリック**により非線形化できるが、結果がカーネルの選び方に依存する