



# データマイニング

## Data Mining

### 12: クラスタリング② Clustering

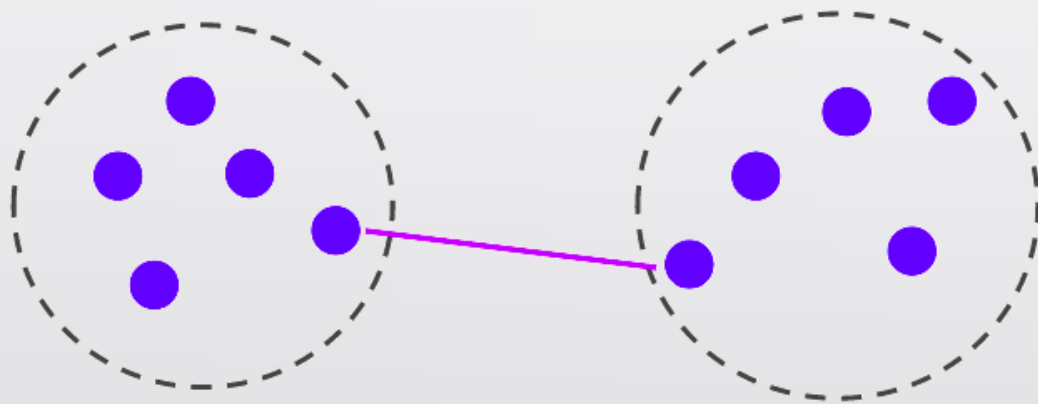
土居 裕和 Hirokazu Doi

長岡技術科学大学 Nagaoka University of Technology

# クラスタ間の距離 Distance between Clusters

## 単リンク法 Single Linkage

Simple linkage



Cluster 1

Cluster 2

$$D(A, B) = \min_{x \in A, y \in B} d(x, y)$$

各クラスタのデータの内、最も近いデータ間の距離を、クラスタ間の距離とする

Distance between clusters is defined as the distance between their closest members

## 単リンク法 Single Linkage

- 大きなクラスターが形成されやすい Large cluster is likely to be formed
- 近いデータ同士が別のクラスターに含まれてしまう連鎖効果が起きやすい  
Neighboring data points tend to be included in separate clusters (chain effect)

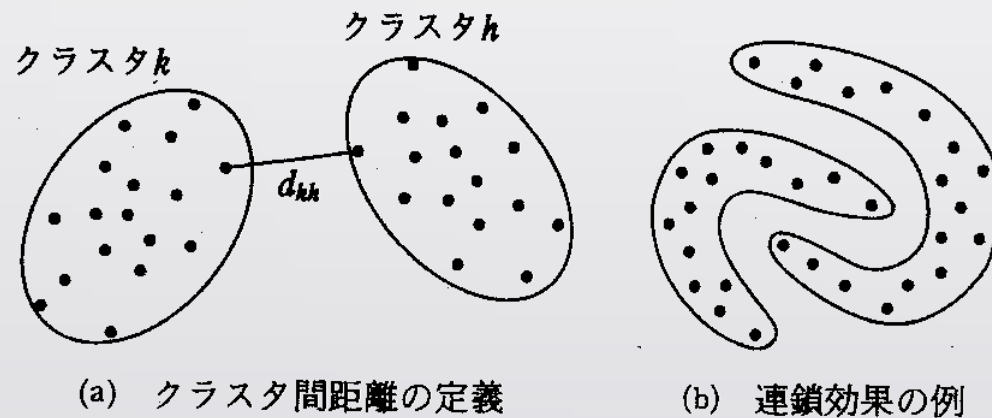
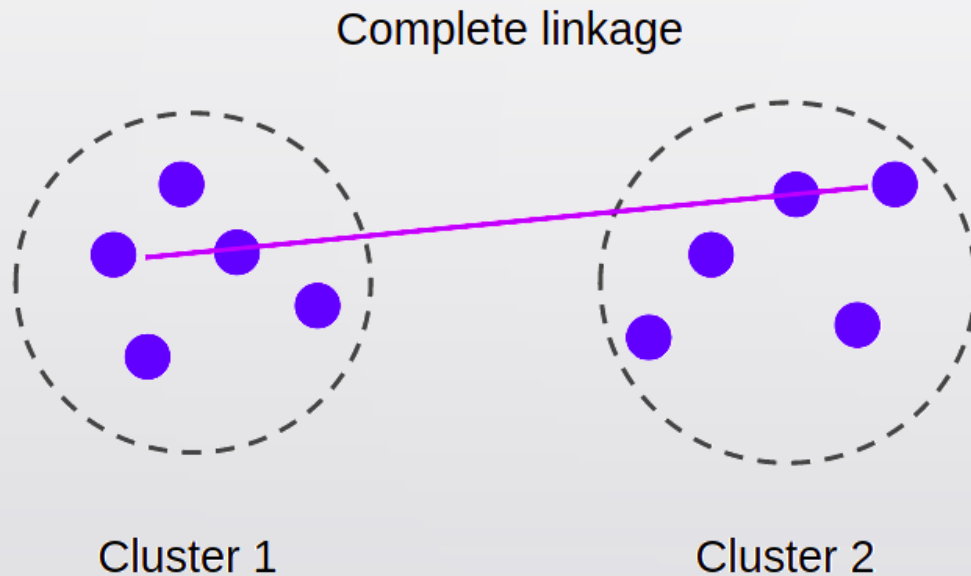


図 1.2.5 最短距離法におけるクラスター間距離の定義と連鎖効果

<https://www.is.kochi-u.ac.jp/kyoko/edu/image/c.html>

# クラスタ間の距離 Distance between Clusters

## 完全リンク法 Complete Linkage



$$D(A, B) = \max_{x \in A, y \in B} d(x, y)$$

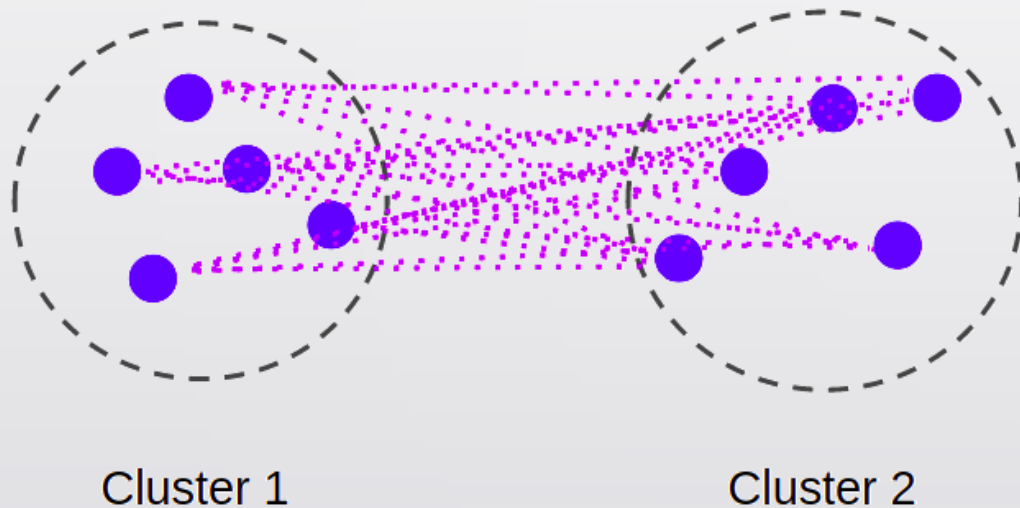
各クラスタのデータの内、最も遠いデータ間の距離を、クラスタ間の距離とする

Distance between clusters is defined as the distance between their farthest members

## クラスター間の距離 Distance between Clusters

### 平均リンク法 Average Linkage

Average linkage



$$D(A, B) = \frac{1}{N_A N_B} \sum_{x \in A} \sum_{y \in B} d(x, y)$$

各クラスターのすべてのデータペアの平均

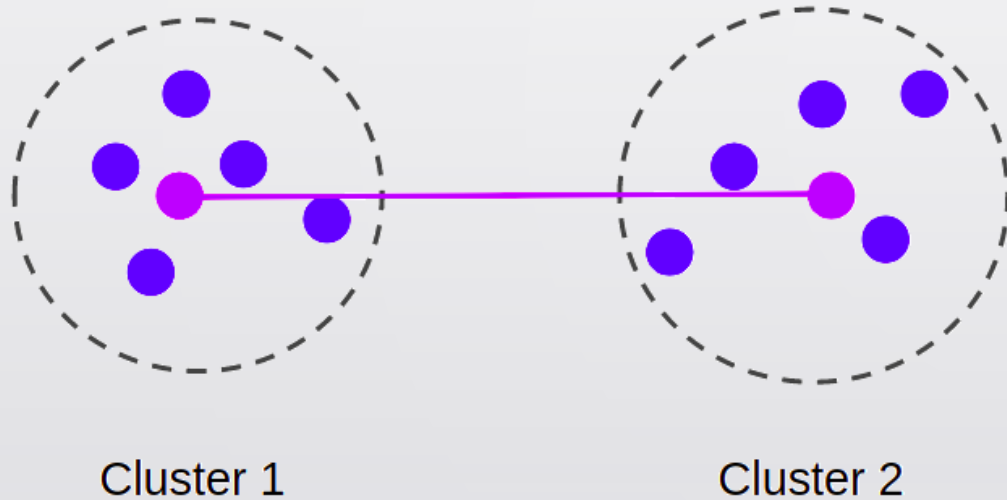
Distance between clusters is defined as average distance of all the between-cluster data pairs



## クラスター間の距離 Distance between Clusters

### 中心リンク法 Centroid Linkage

Centroid linkage



$$D(A, B) = d(\mu_A, \mu_B)$$

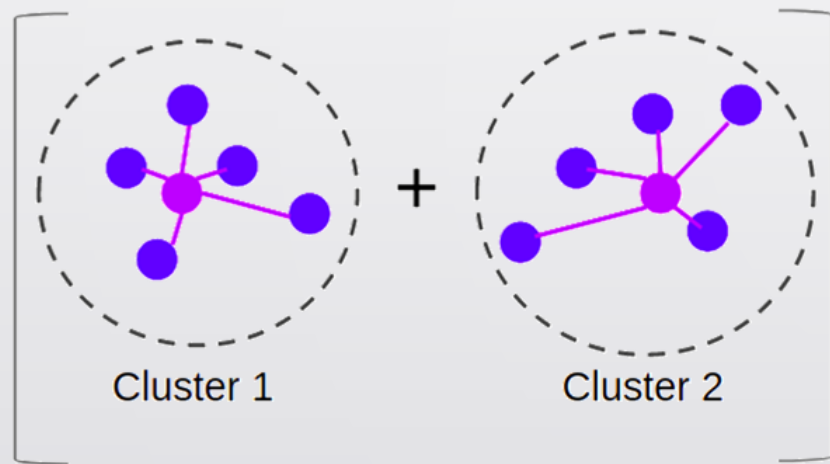
各クラスターの中心間の距離

Distance between clusters is defined as the distance between cluster centers

## ウォード法 Ward Linkage

$\Delta$ が最小になるようなクラスター同士を結合する

Link clusters with minimum  $\Delta$



クラスター内SSEの合計を計算する

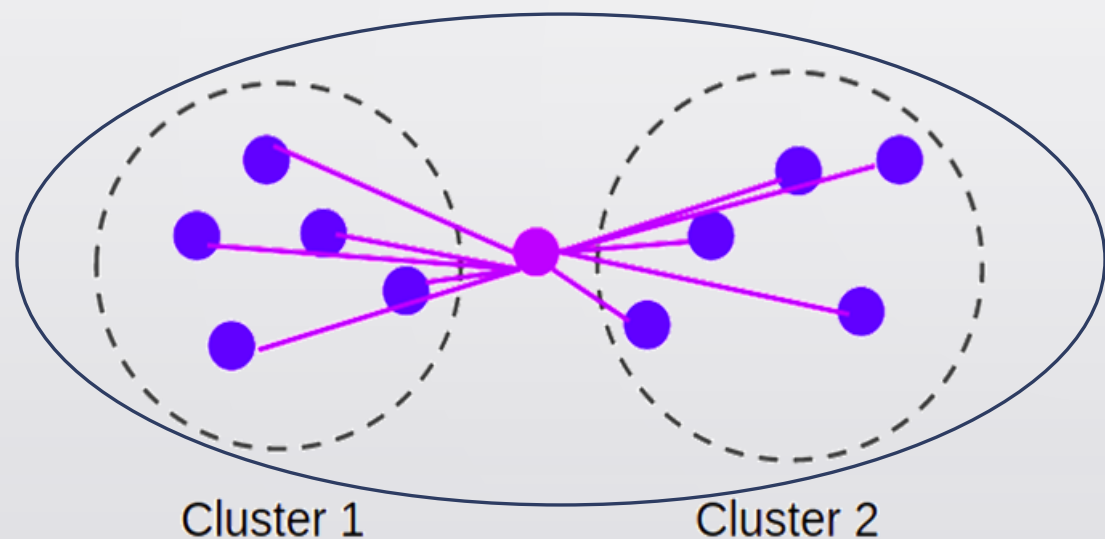
Compute sum of intra-cluster sum of squared error (SSE)

$$\sum_{x \in A} d(x, \mu_A)^2 + \sum_{y \in B} d(y, \mu_B)^2$$

## ワード法 Ward Linkage

$\Delta$ が最小になるようなクラスター同士を結合する

Link clusters with minimum  $\Delta$



クラスターを結合した時のSSEを計算する

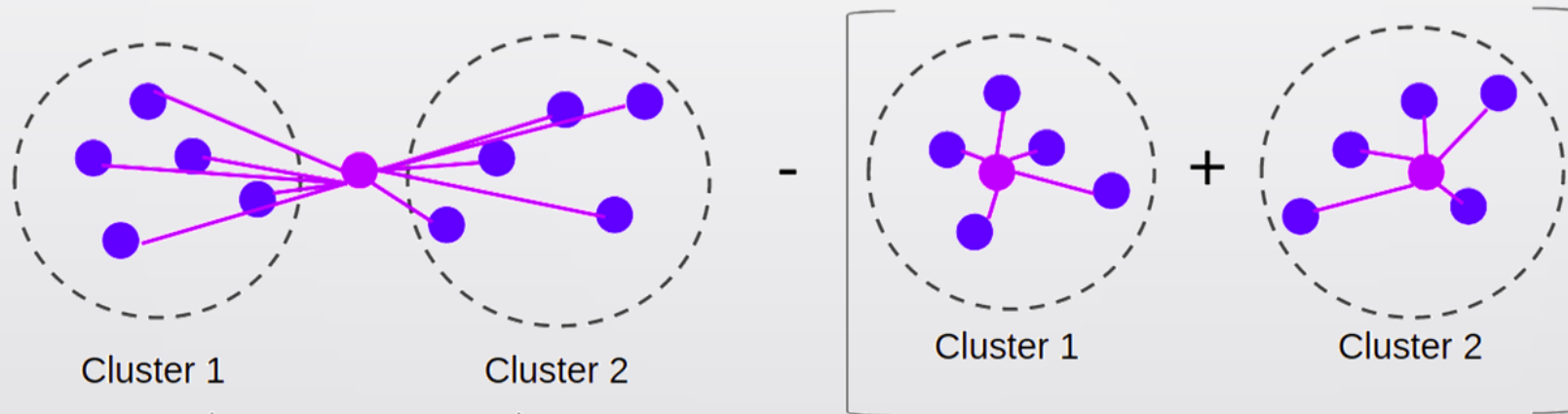
Compute within-cluster sum of squared error (SSE) when the two clusters are joined to form single cluster

$$\sum_{x \in AB} d(x, \mu_{AB})^2$$



## ウォード法 Ward Linkage

$\Delta$ が最小になるようなクラスター同士を結合する  
Link clusters with minimum  $\Delta$

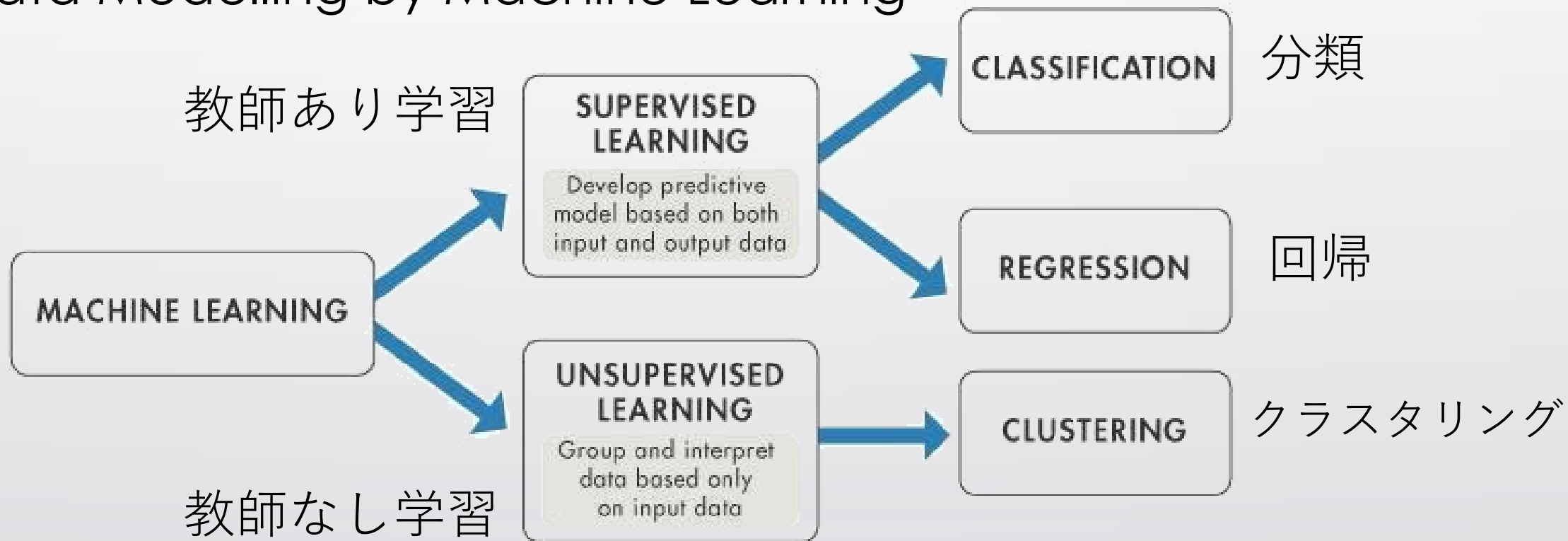


$$\Delta = SSE \text{ after linkage} - SSE \text{ before linkage}$$

$\Delta$ を情報ロスと呼ぶ  $\Delta$  is referred to as "Information loss"

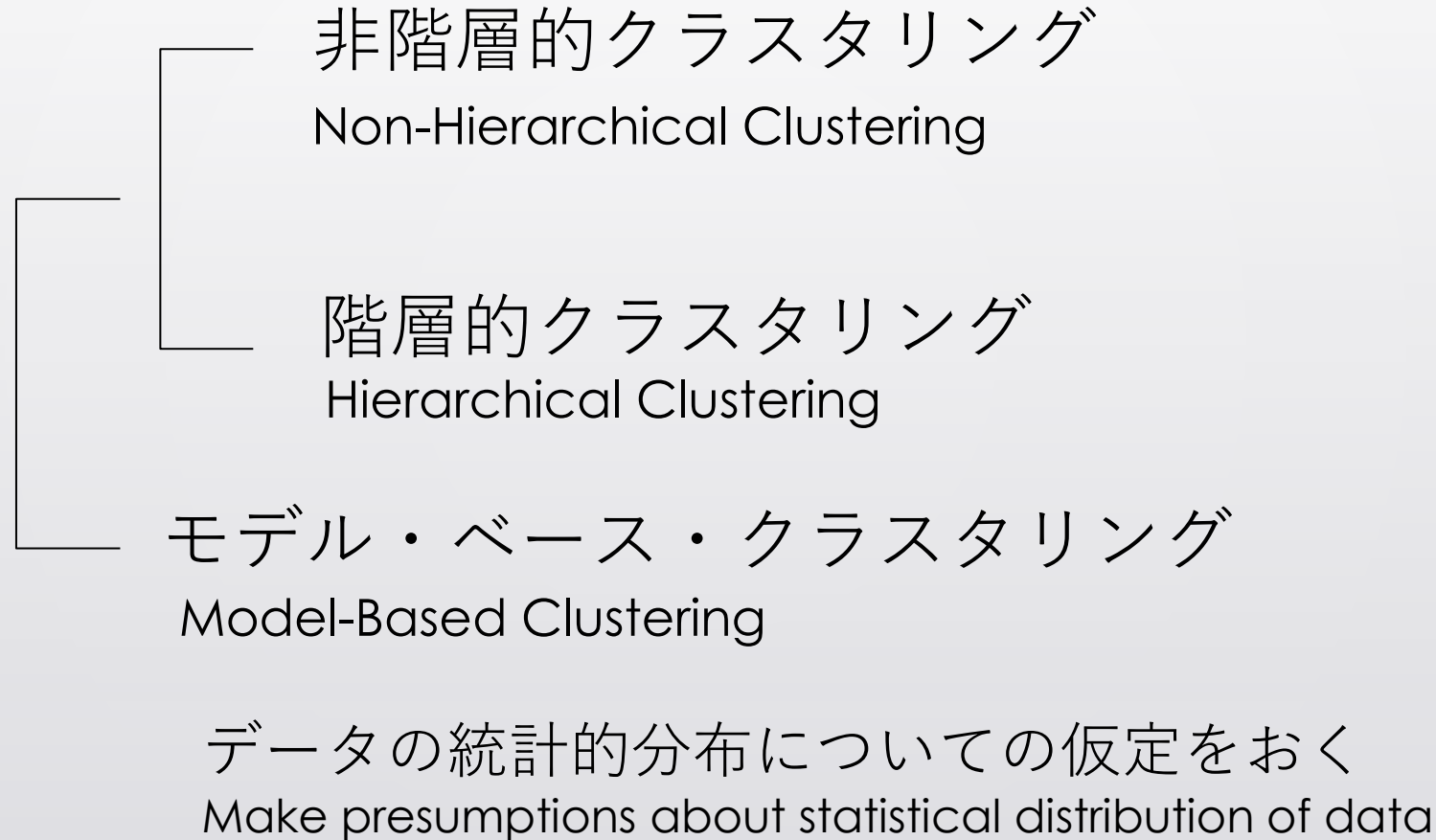
# 機械学習によるモデル化

## Data Modelling by Machine Learning



Supervised Learning versus Unsupervised Learning (Mathworks, n.d.)

# クラスタリングの種類 Types of Clustering





## ソフトクラスタリング Soft Clustering

### ハードクラスタリング Hard Clustering

各データは一つのクラスターにしか所属できない

Each data belongs to single cluster

### ソフトクラスタリング Soft Clustering

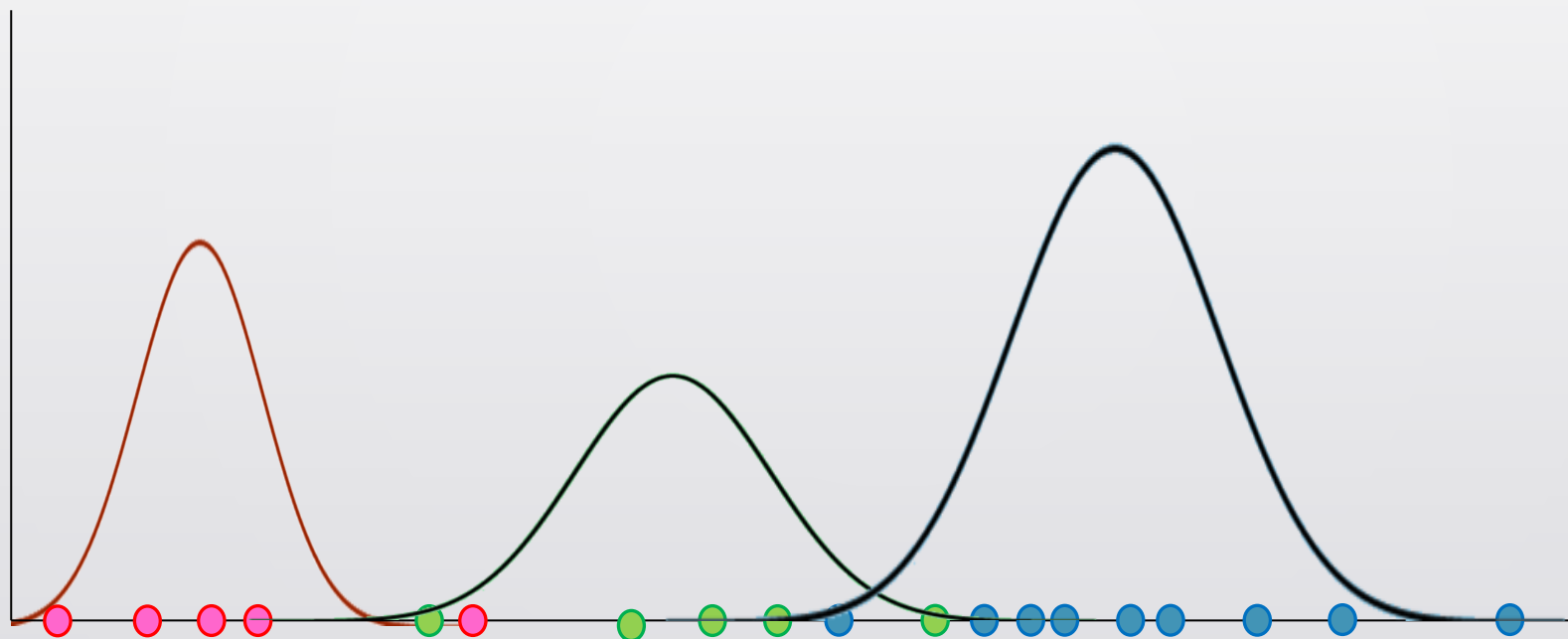
各データが複数のクラスに所属しうる

Each data can belong to multiple classes

# 混合ガウス分布モデル Gaussian Mixture Model

観測されたデータが複数のガウス分布の重ね合わせから生成されたと仮定する

Assume that observations are generated by multiple overlapping Gaussian distributions





# 多次元ガウス分布 Multidimensional Gaussian Distribution

正規分布を多次元に拡張した分布

Probability density distribution obtained by extending normal distribution to multi-dimensional space

$$N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

$\boldsymbol{\Sigma}$  : 分散共分散行列 Variance-covariance matrix

$|\boldsymbol{\Sigma}|$  : 分散共分散行列の行列式 Determinant of variance-covariance matrix

$\boldsymbol{\Sigma}^{-1}$  : 分散共分散行列の逆行列 Inverse matrix of variance-covariance matrix

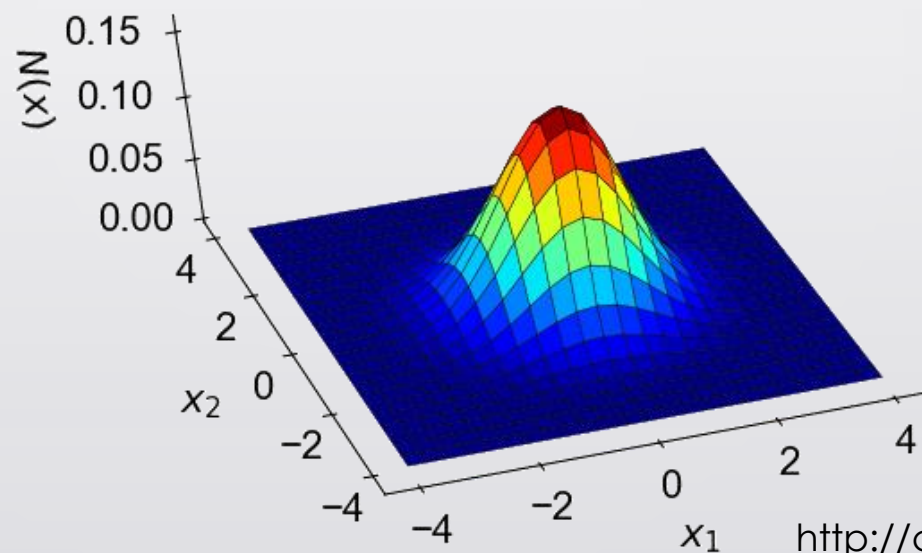
# 多次元ガウス分布 Multidimensional Gaussian Distribution

$$N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{2\pi\sqrt{|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\mathbf{x} = (x_1, x_2)$$

$$\boldsymbol{\mu} = (\mu_1, \mu_2)$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} \\ \sigma_{21} & \sigma_{22}^2 \end{pmatrix}$$



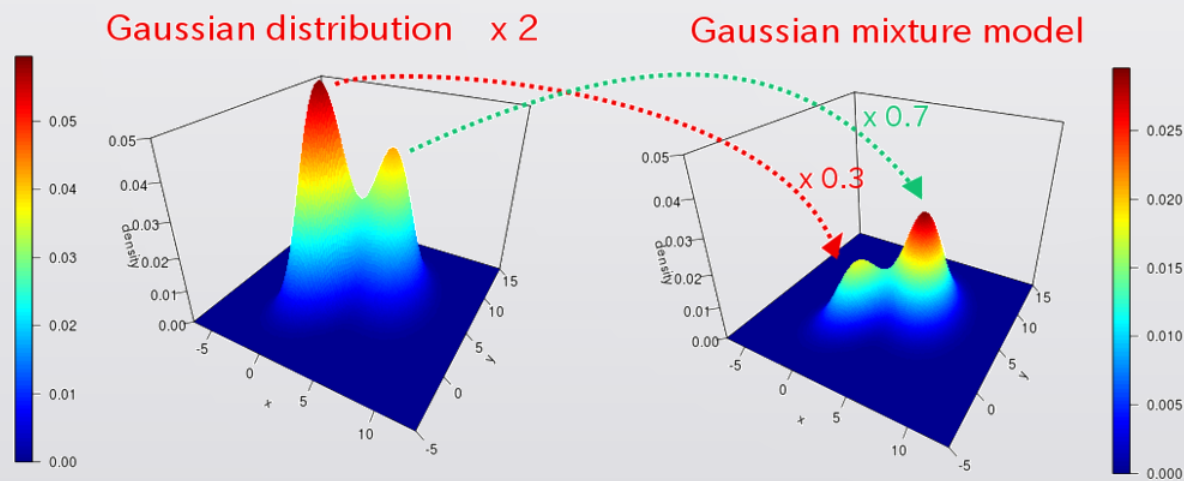
<http://chasen.org/~daiti-m/gpbook/>

# 混合ガウス分布 Gaussian Mixture Distribution

$M$ 個の正規分布の重ね合わせにより確率分布を表現する

Represent probability distribution as weighted mixture of  $M$  normal distributions

$$p(\mathbf{x}) = \sum_{m=1}^M \pi_m N(\mathbf{x} | \boldsymbol{\mu}_m, \boldsymbol{\sigma}_m) \quad 0 \leq \pi_m \leq 1 \quad \sum_{m=1}^M \pi_m = 1 \quad \pi_m : \text{混合比 Mixing Ratio}$$



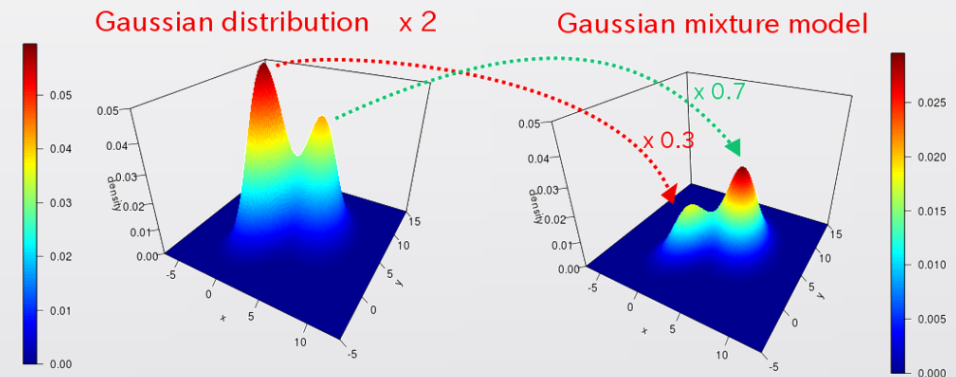
<https://work-in-progress.hatenablog.com/entry/2018/11/08/224826>

# 潜在変数 Latent Variable (隠れ変数 Hidden Variable)

観測データからは直接得ることが出来ない情報

Information that cannot be obtained directly from observations

$$p(\mathbf{x}) = \sum_{m=1} \pi_m N(\mathbf{x} | \boldsymbol{\mu}_m, \boldsymbol{\sigma}_m)$$



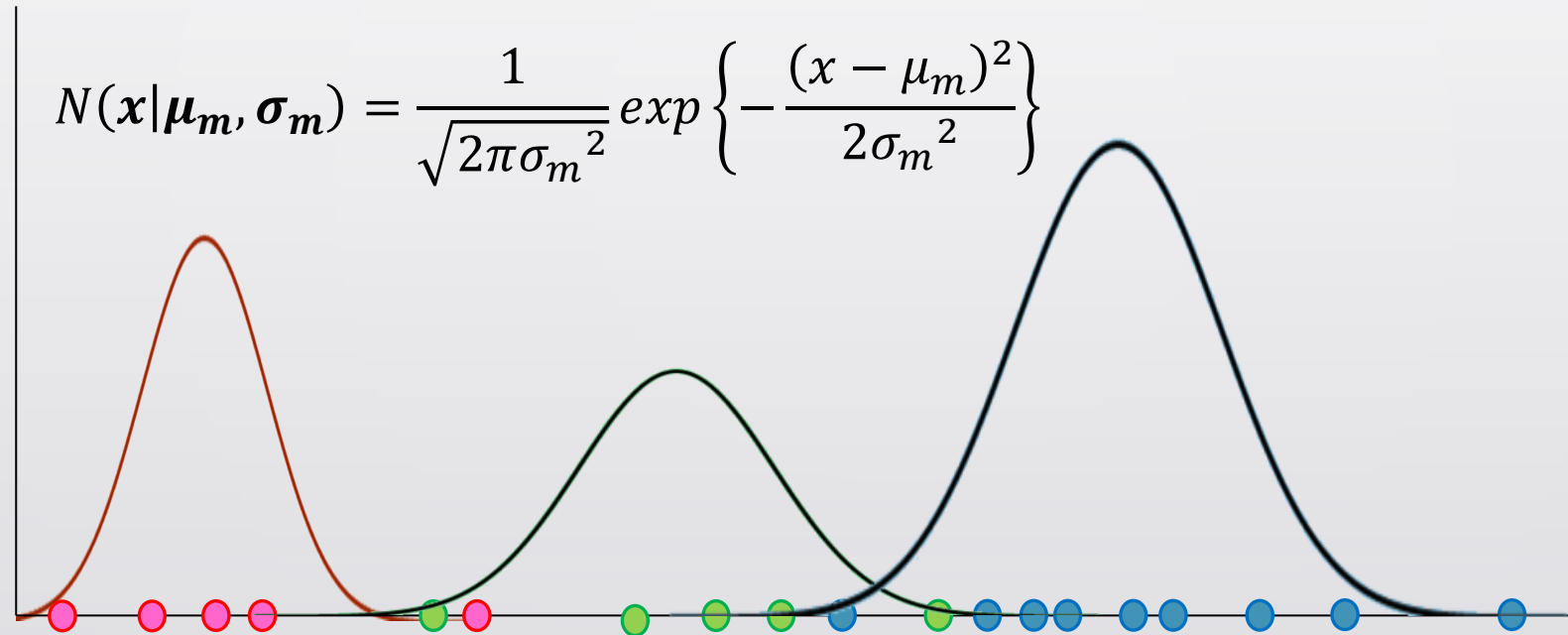
潜在変数を使って、どの分布からデータが生成されたかを表現する

Latent variable represents from which gaussian distribution a data is generated¥

# 一次元混合ガウスモデル 1-Dimensional GMM

観測されたデータが複数のガウス分布の重ね合わせから生成されたと仮定する

Assume that observations are generated by multiple overlapping Gaussian distributions





## 一次元混合ガウスモデル 1-Dimensional GMM

$\pi_m$  : 混合比 Mixing Ratio  $\sum_{m=1}^M \pi_m = 1 \quad 0 \leq \pi_m \leq 1$

$\mathbf{z}$ : 潜在変数 Latent Variables

$$z_m \in \{0, 1\} \quad \mathbf{z} = (z_1, z_2, z_3, z_4, \dots, z_{M-1}, z_M)$$

$$\sum_{m=1}^M z_m = 1 \quad \text{ex) } \mathbf{z} = (0, 0, 0, 1, \dots, 0, 0, 0)$$

## 一次元混合ガウスモデル 1-Dimensional GMM

$$N(x|\mu_m, \sigma_m) = \frac{1}{\sqrt{2\pi\sigma_m^2}} \exp\left\{-\frac{(x - \mu_m)^2}{2\sigma_m^2}\right\}$$

$$p(x) = \sum_{m=1}^M p(x|z_m = 1)p(z_m = 1) = \sum_{m=1}^M N(x|\mu_m, \sigma_m)\pi_m$$

## 一次元混合ガウスモデル 1-Dimensional GMM

$$p(z_m = 1|x) = \frac{p(x|z_m = 1)p(z_m = 1)}{p(x)} = \frac{p(x|z_m = 1)p(z_m = 1)}{\sum_{m=1}^M p(x|z_m = 1)\pi_m}$$

観測された $x$ が分布 $m$ から生成された事後確率

Posterior probability that observation  $x$  is generated by distribution  $m$

潜在変数の期待値は、その事後確率と一致する

Expected value of latent variable corresponds to its posterior probability

$$E[z_m] = p(z_m = 1|x)$$

## 完全データ Complete Data

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2 \cdots \mathbf{x}_N\} \quad \mathbf{Z} = \{\mathbf{z}_1, \mathbf{z}_2 \cdots \mathbf{z}_N\} \quad \mathbf{z}_n = (z_{n1}, z_{n2} \cdots z_{nM})$$

$$\mathbf{Y} = \{\mathbf{X}, \mathbf{Z}\} = \{\mathbf{x}_1, \mathbf{x}_2 \cdots \mathbf{x}_N, \mathbf{z}_1, \mathbf{z}_2 \cdots \mathbf{z}_N\}$$

$$p(\mathbf{x}_n, z_{nm} = 1 | \boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi})$$

$$= p(\mathbf{x}_n | z_{nm} = 1) p(z_{nm} = 1 | \boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi})$$

$$= N(\mathbf{x}_n | \boldsymbol{\mu}_m, \boldsymbol{\sigma}_m) \pi_m$$

## MLEによるパラメータ推定Parameter Estimation by MLE

$$p(\mathbf{x}_n, z_{nm} = 1 | \boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi}) = N(\mathbf{x}_n | \mu_m, \sigma_m) \pi_m$$

$$p(\mathbf{Y} | \boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi}) = \prod_{n=1}^N \prod_{m=1}^M [N(\mathbf{x}_n | \mu_m, \sigma_m) \pi_m]^{z_{nm}}$$

$$\boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi} = \operatorname{argmax}_{\boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi}} p(\mathbf{Y} | \boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi})$$

観測される $\mathbf{X}$ と対応する $\mathbf{Z}$ の同時確率を最大化するパラメータセットを求める

Search for parameter set that maximizes observations  $\mathbf{X}$  and corresponding  $\mathbf{Z}$



## MLEによるパラメータ推定Parameter Estimation by MLE

$$\log p(Y | \boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{m=1}^M \mathbf{z}_{nm} \log N(\mathbf{x}_n | \mu_m, \sigma_m) + \sum_{n=1}^N \sum_{m=1}^M \mathbf{z}_{nm} \log \pi_m$$

潜在変数は直接的に観察できない

Latent variables are not directly observable



潜在変数の期待値 = 事後確率で置き換えてパラメータ推定

Estimate parameters by replacing latent variables with their expected values

//////  
Q関数 Q function

$$\log p(Y | \boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{m=1}^M \mathbf{z}_{nm} \log N(\mathbf{x}_n | \mu_m, \sigma_m) + \sum_{n=1}^N \sum_{m=1}^M \mathbf{z}_{nm} \log \pi_m$$

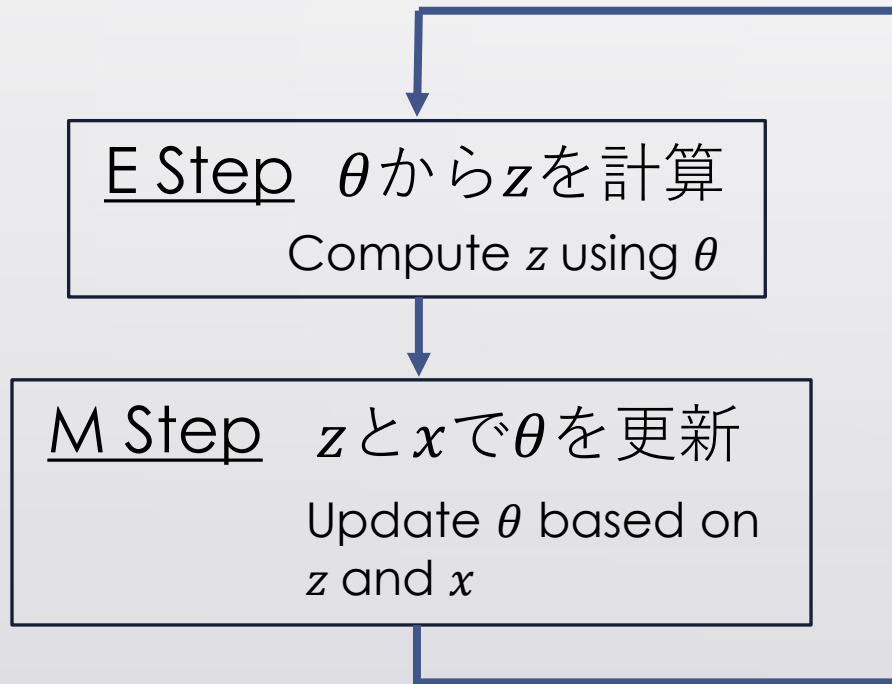


$$Q = \sum_{n=1}^N \sum_{m=1}^M E[\mathbf{z}_{nm}] \log N(\mathbf{x}_n | \mu_m, \sigma_m) + \sum_{n=1}^N \sum_{m=1}^M E[\mathbf{z}_{nm}] \log \pi_m$$

# EM アルゴリズム Expectation-Maximizing Algorithm

潜在変数を含むモデルの代表的なパラメータ推定法

Algorithm for parameter estimation of models including latent variables



$x$ : 観測 Observations

$\theta$ : 確率密度関数のパラメータセット  
Parameter set of probability distribution functions

$z$ : 潜在変数 Latent Variables

## E-ステップ Expectation-Step

現在のパラメータセット  $\theta^{(t)}$  を用いて  $z$  の期待値を求める

Compute expected value of  $z$  based on current parameter set  $\theta^{(t)}$

$$\theta^{(t)} = \{\mu^{(t)}, \sigma^{(t)}, \pi^{(t)}\}$$

$$\mathbf{z} = (z_1, z_2 \cdots z_m \cdots z_M)$$

$$E[z_m] = \frac{N(x | \mu_m^{(t)}, \sigma_m^{(t)}) \pi_m^{(t)}}{\sum_{m=1}^M N(x | \mu_m^{(t)}, \sigma_m^{(t)}) \pi_m^{(t)}}$$

## M-ステップ Maximization-Step

$$\begin{aligned} p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}^{(t)}) &= \prod_{n=1}^N \prod_{m=1}^M p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}^{(t)}) p(\mathbf{z} | \boldsymbol{\theta}^{(t)}) \\ &= \prod_{n=1}^N \prod_{m=1}^M [p(x_n | z_{n,m}, \boldsymbol{\theta}^{(t)}) p(z_m | \boldsymbol{\theta}^{(t)})]^{z_{n,m}} \\ &= \prod_{n=1}^N \prod_{m=1}^M [N(x_n | \mu_m^{(t)}, \sigma_m^{(t)}) \pi_m^{(t)}]^{z_{n,m}} \end{aligned}$$



## M-ステップ Maximization-Step

Q関数を最大化するようパラメータセット  $\theta^{(t)}$  を更新

Update parameter set  $\theta^{(t)}$  so that Q function is maximized

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}^{(t)}) = \prod_{n=1}^N \prod_{m=1}^M [N(x_n | \mu_m^{(t)}, \sigma_m^{(t)}) \pi_m^{(t)}]^{z_{n,m}}$$

$$Q(\boldsymbol{\theta}^{(t)}) = \sum_{n=1}^N \sum_{m=1}^M E[z_{n,m}] \left[ \log(\pi_m^{(t)}) + \log(N(x_n | \mu_m^{(t)}, \sigma_m^{(t)})) \right]$$

## M-ステップ Maximization-Step

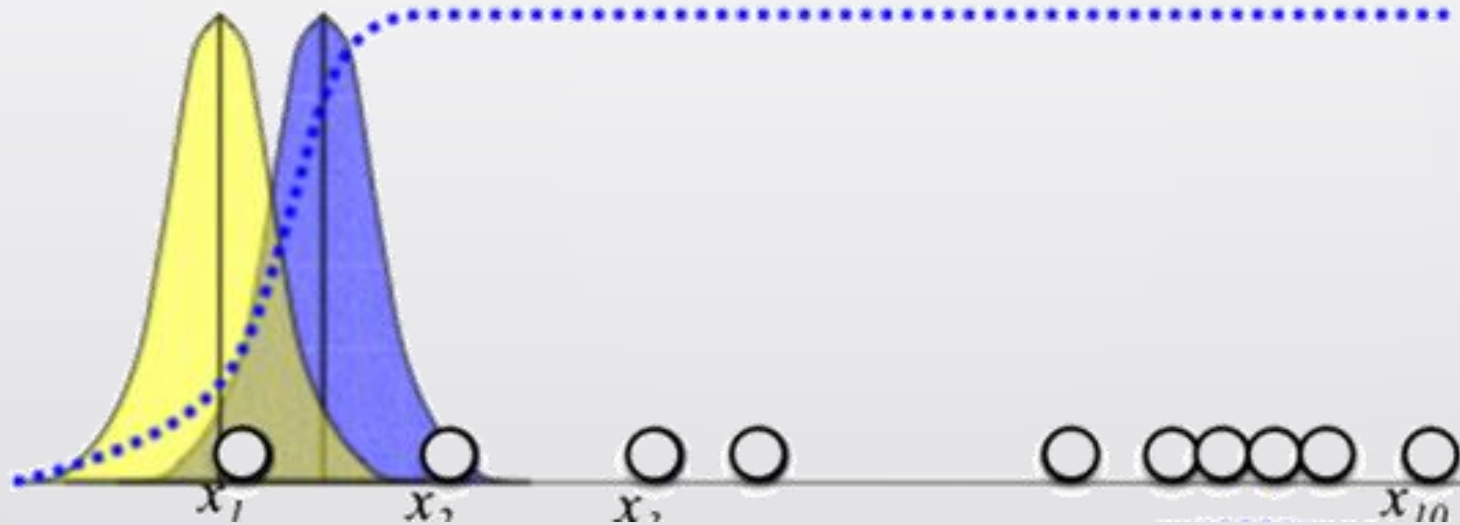
$$Q(\boldsymbol{\theta}^{(t)}) = \sum_{n=1}^N \sum_{m=1}^M E[z_{n,m}] \left[ \log(\pi_m^{(t)}) + \log \left( N(x_n | \mu_m^{(t)}, \sigma_m^{(t)}) \right) \right]$$

$$E[z_m] = p(z_m = 1 | x) = \frac{p(x | z_m = 1) p(z_m = 1)}{\sum_{m=1}^M p(x | z_m = 1) \pi_m}$$

$$\frac{\partial Q}{\partial \pi_m} = 0 \quad \frac{\partial Q}{\partial \mu_m} = 0 \quad \frac{\partial Q}{\partial \sigma_m} = 0$$

# EM アルゴリズム Expectation-Maximizing Algorithm

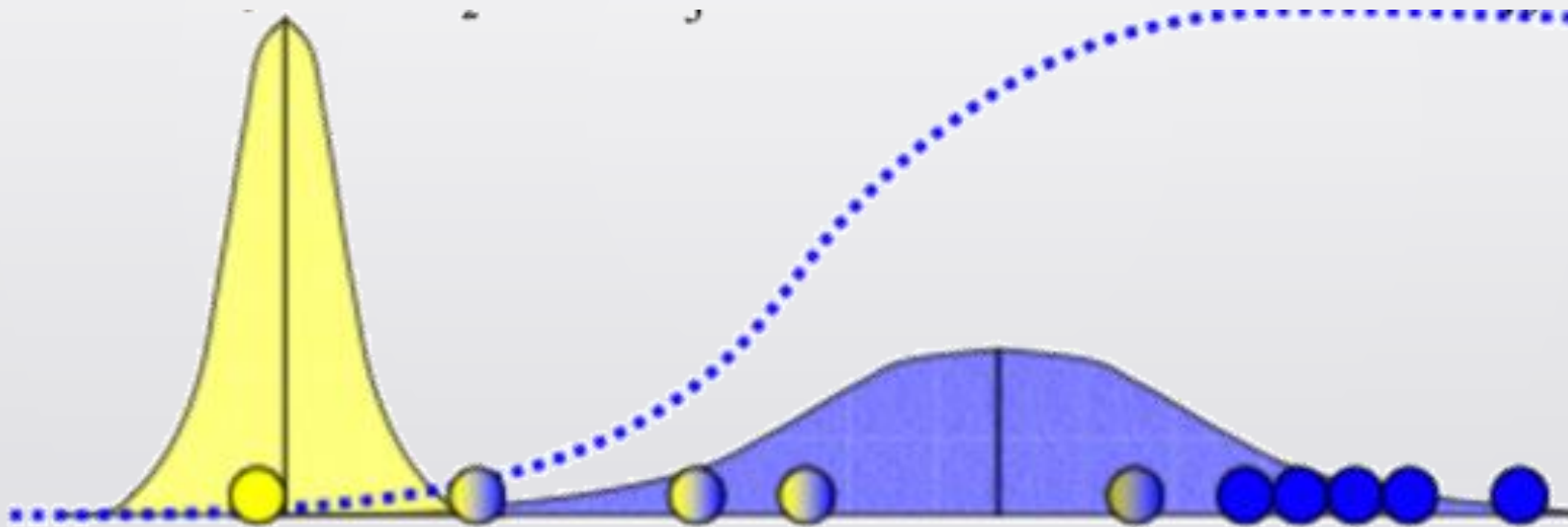
1.  $\theta$  を初期化する Initialize  $\theta$



<https://courses.cs.washington.edu/courses/cse416/22sp/lectures/12/12.pdf>

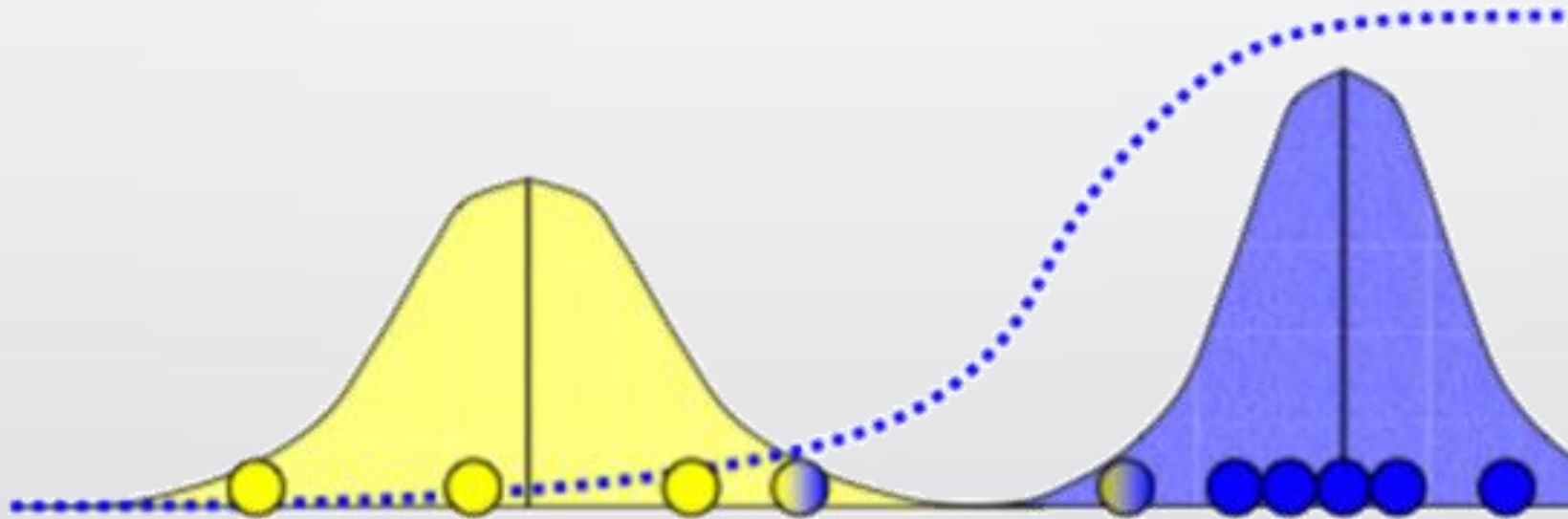
# EMアルゴリズム Expectation-Maximizing Algorithm

2. E-ステップ：現在のパラメータセット $\theta^{(t)}$ を用いて $z$ の期待値を求める  
データの色が $E[z_{n,m}]$ を表す



# EMアルゴリズム Expectation-Maximizing Algorithm

3. M-ステップ：パラメータセット  $\theta^{(t)}$  を更新する    Update  $\theta$



<https://courses.cs.washington.edu/courses/cse416/22sp/lectures/12/12.pdf>



## クラスタリングの評価 Evaluation of Clustering Results

	$c^1$	$c^2$
a	90	10
b	20	80

$C^k$ :  $k$ 番目のクラスター  $k$ -th cluster

$n_{m,k}$ :  $C^k$ のうち $m$ 番目のクラス $C_m$ に属するデータの数  
Number of data belonging to class  $C_m$  within  $C^k$

$|C^k|$ : クラスター $C^k$ に属するデータの数  
Number of data belonging to cluster  $C^k$

## 純度 Purity

### 局所的純度 Local Purity

$$Purity = \frac{\max_m n_{m,k}}{|C^k|}$$

最大多数派のクラスがクラスターに占める割合

Proportion of data of majority class

### 大域的純度 Global Purity

$$Purity = \frac{\sum_k \max_m n_{m,k}}{\sum_k |C^k|}$$

	$C^1$	$C^2$
a	90	10
b	20	80

## 逆純度 Inverse Purity

クラスターの純度は2つのテーブルで同じ

Purity of the clusters is the same across the tables below

	$c^1$	$c^2$
a	90	10
b	20	80

	$c^1$	$c^2$
a	90	40
b	20	5

$$Purity = \frac{\max_m n_{m,k}}{|C^k|}$$

## 逆純度 Inverse Purity

$M_k = \operatorname{argmax}_m n_{m,k}$  クラスタ  $C^k$  において最大多数派のクラス  
Majority class within cluster  $C_k$

$|C_{M_k}|$ : クラス  $M_k$  に属するデータの総数  
Total number of data belonging to class  $M_k$

$$|C_{M_k}| = \sum_{k=1}^K n_{M_k,k}$$

## 逆純度 Inverse Purity

$$\text{Inverse Purity} = \frac{1}{N} \sum_{k=1}^K \frac{\max_m n_{m,k}}{|C_{M_k}|} |C^k|$$

$M_1 = a$

	$c^1$	$c^2$
a	90	10
b	20	80

$M_2 = b$

	$c^1$	$c^2$
a	90	10
b	20	80



## 逆純度 Inverse Purity

$$\text{Inverse Purity} = \frac{1}{N} \sum_{k=1}^K \frac{\max_m n_{m,k}}{|C_{M_k}|} |C^k|$$

$M_1 = a$

	$c^1$	$c^2$
a	90	40
b	20	5

$M_2 = a$

	$c^1$	$c^2$
a	90	40
b	20	5

## F值 *F*-value

$$Purity = \frac{\sum_k \max_m n_{m,k}}{\sum_k |C^k|} \quad Inverse Purity = \frac{1}{N} \sum_{k=1}^K \frac{\max_m n_{m,k}}{|C_{M_k}|} |C^k|$$

$$F = \frac{1}{\frac{1}{2} \left( \frac{1}{Purity} + \frac{1}{Inverse Purity} \right)} = \frac{2 Purity \cdot Inverse Purity}{Purity + Inverse Purity}$$