



データマイニング

Data Mining

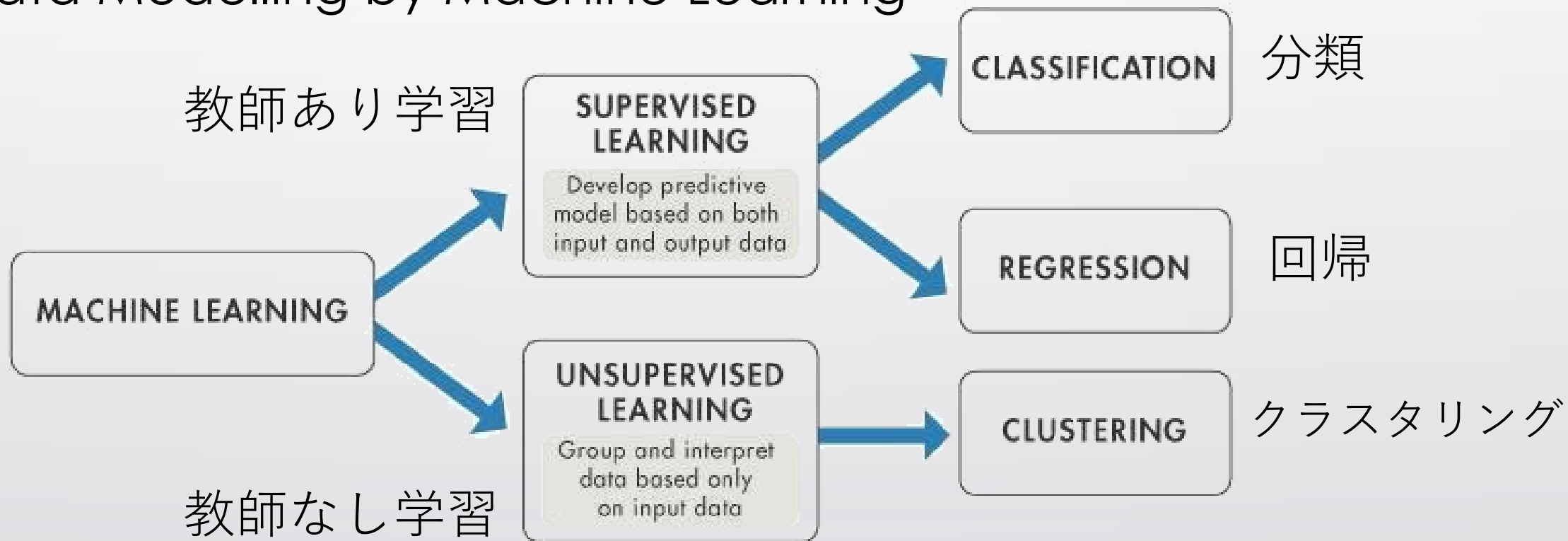
11: クラスタリング① Clustering

土居 裕和 Hirokazu Doi

長岡技術科学大学 Nagaoka University of Technology

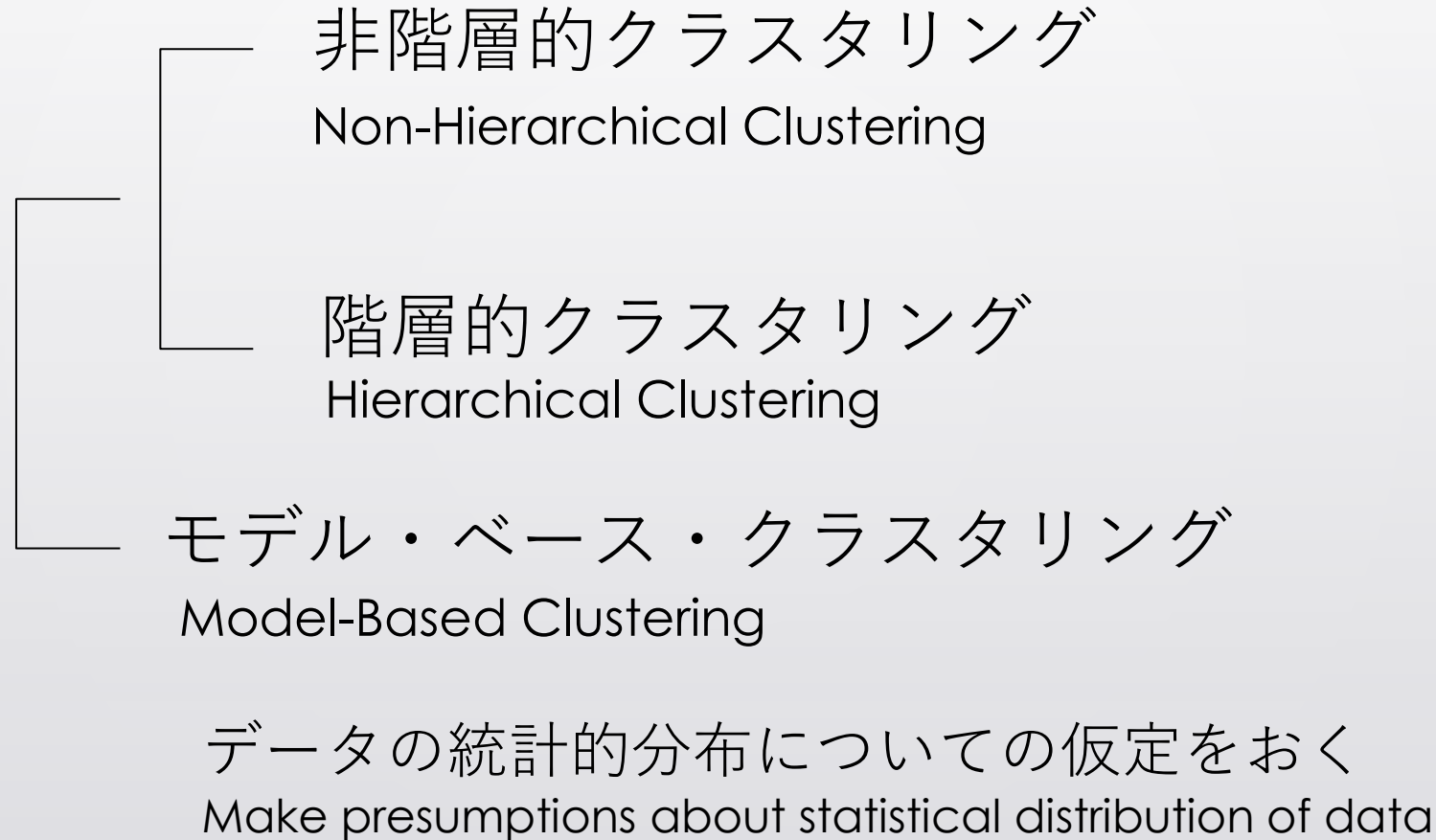
機械学習によるモデル化

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





	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
...
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	virginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica



	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2
...
145	6.7	3.0	5.2	2.3
146	6.3	2.5	5.0	1.9
147	6.5	3.0	5.2	2.0
148	6.2	3.4	5.4	2.3
149	5.9	3.0	5.1	1.8

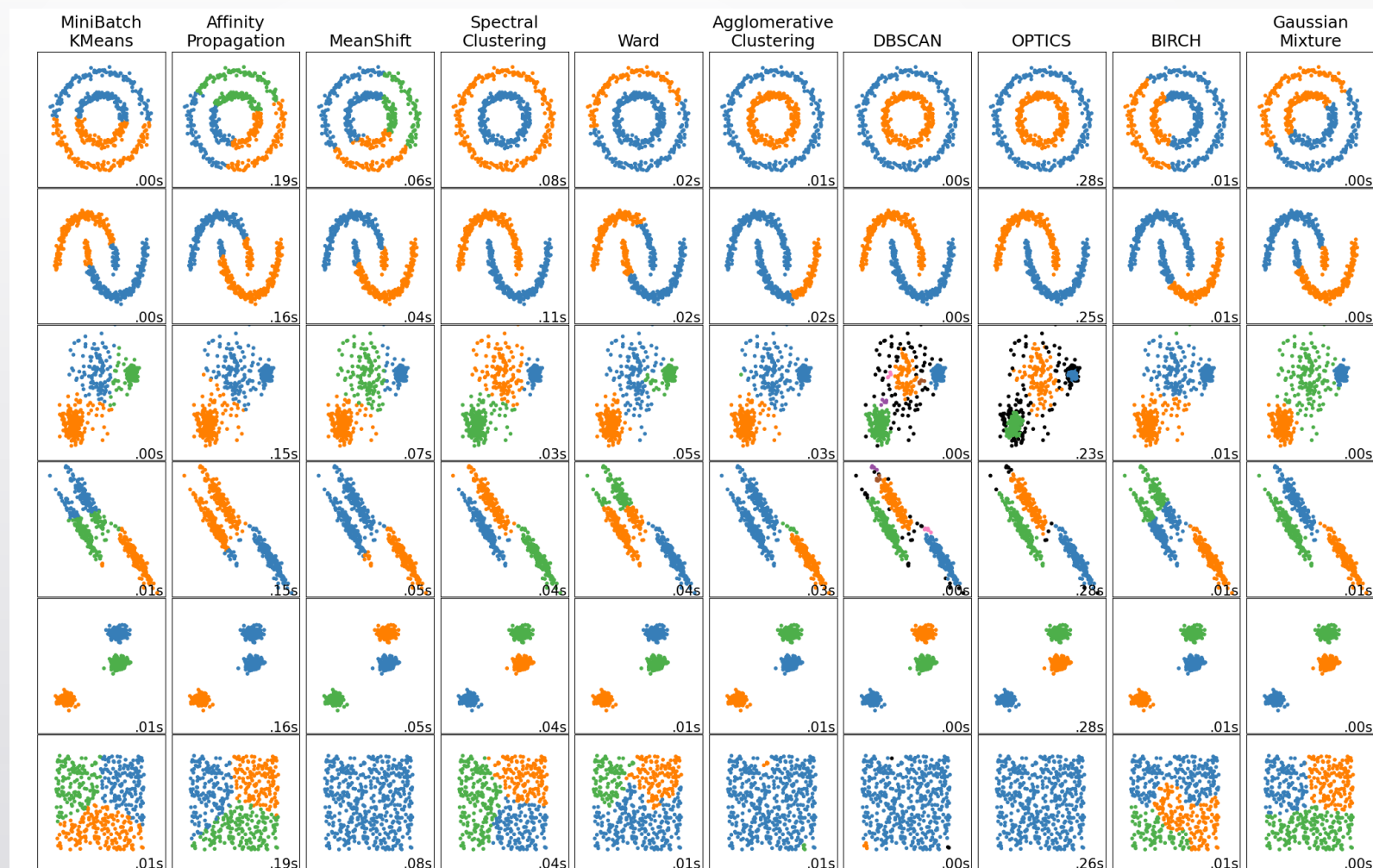
クラスタリング

Final result of clustering
depends on

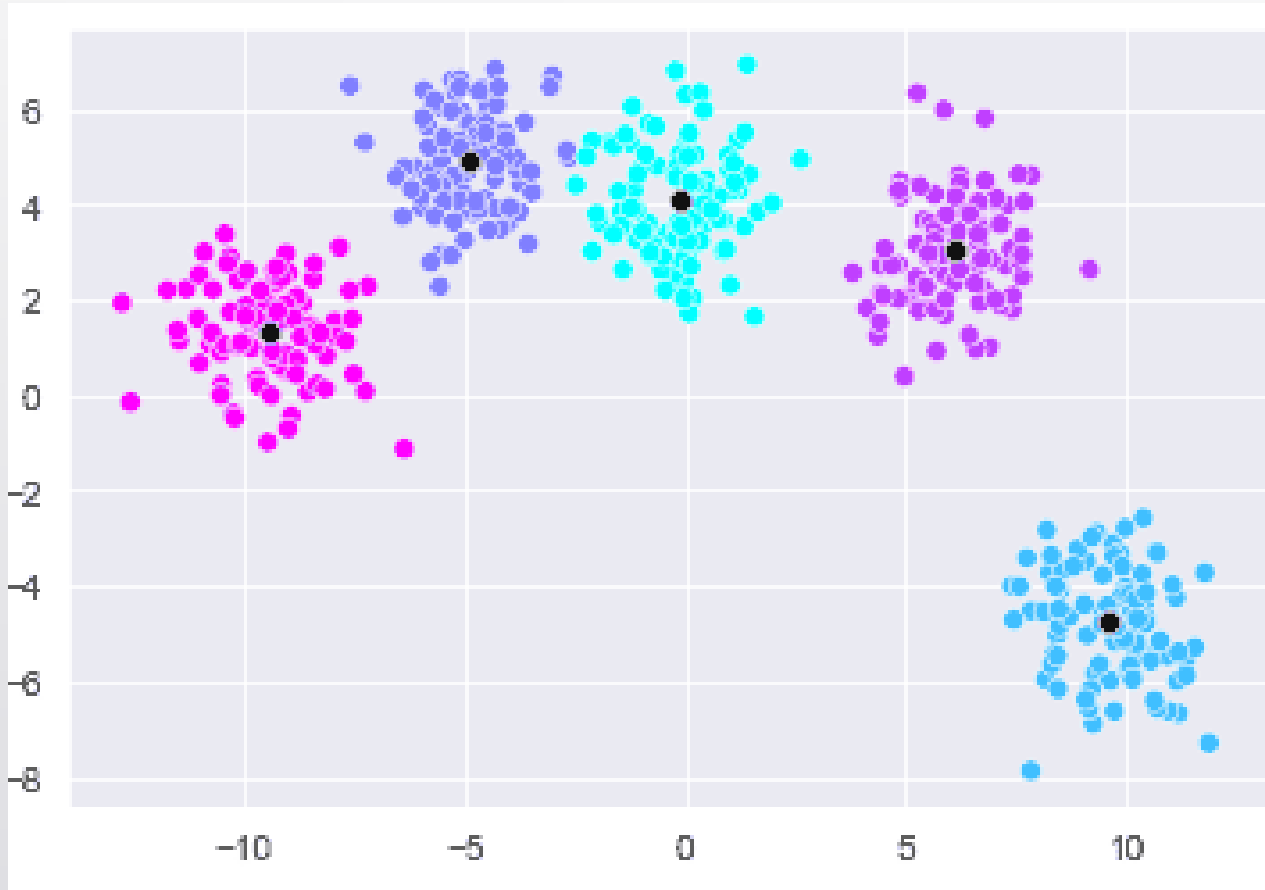
Type of algorithm
アルゴリズムの種類

Parameter Setting
パラメータ設定

<https://scikit-learn.org/stable/modules/clustering.html>




K平均クラスタリング k-means clustering



クラスターの数指定しなくては
いけない

You have to specify the number of
clusters, k .



K平均クラスタリング k-means clustering

非階層的クラスタリングの代表的なアルゴリズム

Representative algorithm of non-hierarchical clustering

各クラスターの中心とデータとの距離に基づいてクラスタリングを行う

Clustering based on distance between data point and center of each cluster

予めクラスターの数进行指定する必要がある

It is necessary to specify the number of clusters beforehand

K平均クラスタリング *k*-means clustering

$$D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}, \mathbf{x}_i \in R^d$$

d次元データがN個ある There are N d-dimensional data points

$\boldsymbol{\mu}_k$: k 番目のクラスターの代表ベクトル
Representative vector of k -th cluster

$M(\boldsymbol{\mu}_k)$: $\boldsymbol{\mu}_k$ のボロノイ領域
Voronoi region of representative vector $\boldsymbol{\mu}_k$

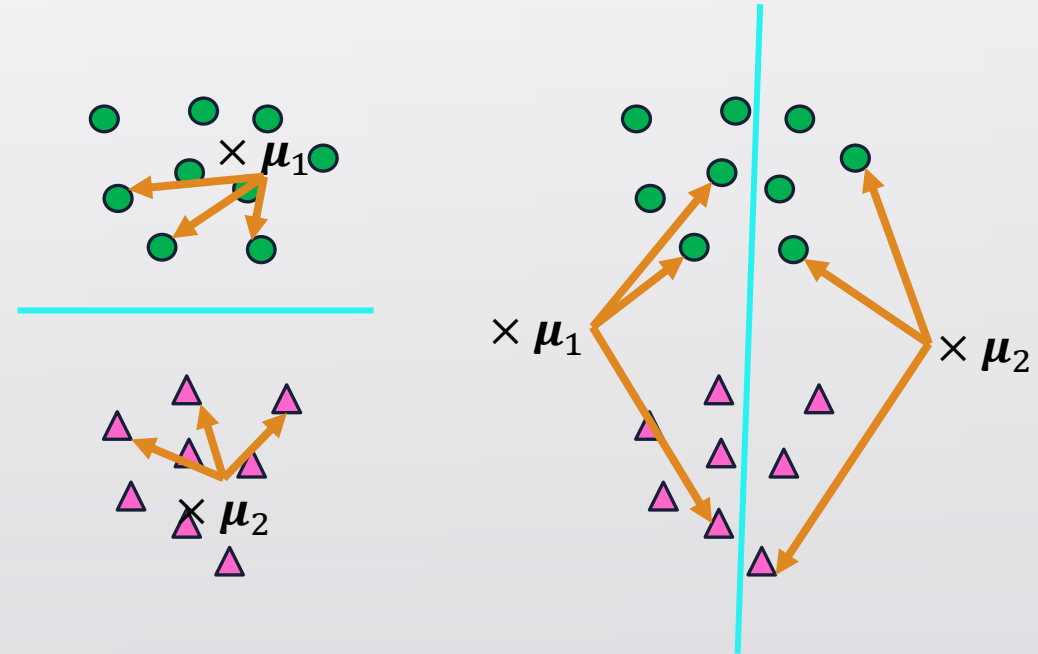
K平均クラスタリング k-means clustering

$$q_{i,k} = \begin{cases} 1 & (x_i \in M(\boldsymbol{\mu}_k) \text{ の場合}) \\ 0 & \text{In case of } x_i \in M(\boldsymbol{\mu}_k) \end{cases}$$

$$J(q_{i,k}, \boldsymbol{\mu}_k) = \sum_{i=1}^N \sum_{k=1}^K q_{i,k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2$$

$J(q_{i,k}, \boldsymbol{\mu}_k)$ を最小化する $q_{i,k}$ と $\boldsymbol{\mu}_k$ を求める

Find $q_{i,k}$ and $\boldsymbol{\mu}_k$ that minimize $J(q_{i,k}, \boldsymbol{\mu}_k)$



K平均クラスタリング k-means clustering

$$J(q_{i,k}, \mu_k) = \sum_{i=1}^N \sum_{k=1}^K q_{i,k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2$$

$$\frac{\partial J}{\partial \boldsymbol{\mu}_k} = 0 \quad -2 \sum_{i=1}^N q_{i,k} (\mathbf{x}_i - \boldsymbol{\mu}_k) = 0$$

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^N q_{i,k} \mathbf{x}_i}{\sum_{i=1}^N q_{i,k}}$$

$\boldsymbol{\mu}_k$ と $q_{i,k}$ を同時に最適化するにはどうすればいいか？

How can we optimize $\boldsymbol{\mu}_k$ and $q_{i,k}$ simultaneously?

K平均クラスタリング k-means clustering

$$J(q_{i,k}, \mu_k) = \sum_{i=1}^N \sum_{k=1}^K q_{i,k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2$$

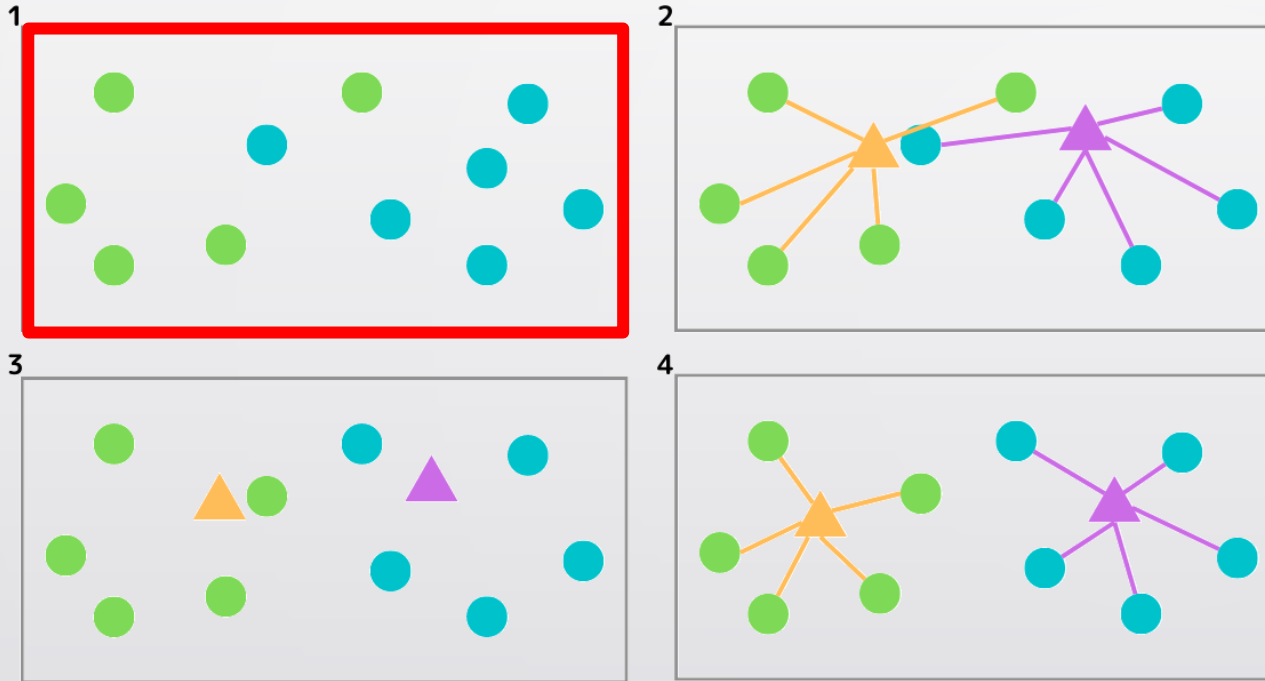
$$\frac{\partial J}{\partial \boldsymbol{\mu}_k} = 0 \quad -2 \sum_{i=1}^N q_{i,k} (\mathbf{x}_i - \boldsymbol{\mu}_k) = 0$$

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^N q_{i,k} \mathbf{x}_i}{\sum_{i=1}^N q_{i,k}}$$

$\boldsymbol{\mu}_k$ と $q_{i,k}$ を同時に最適化するにはどうすればいいか？

How can we optimize $\boldsymbol{\mu}_k$ and $q_{i,k}$ simultaneously?

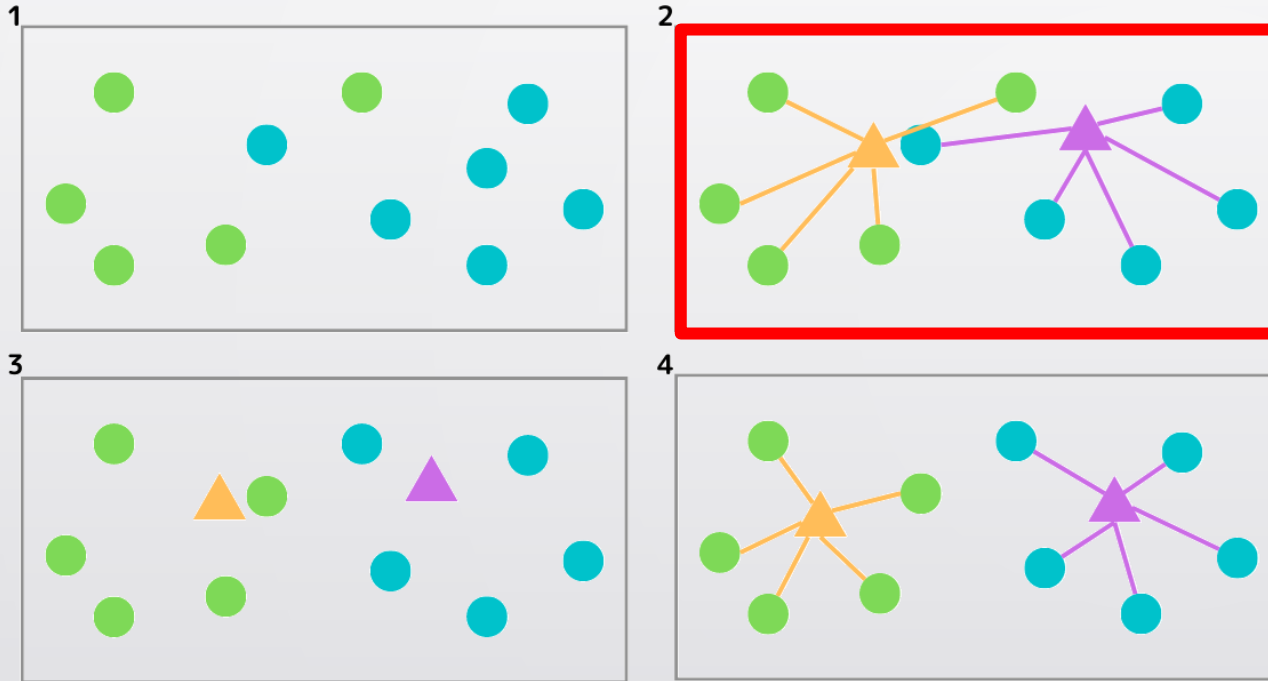
K平均クラスタリング k-means clustering



1. データをランダムにクラスターに割り当てる

Randomly assign data to one of the clusters

K平均クラスタリング k-means clustering

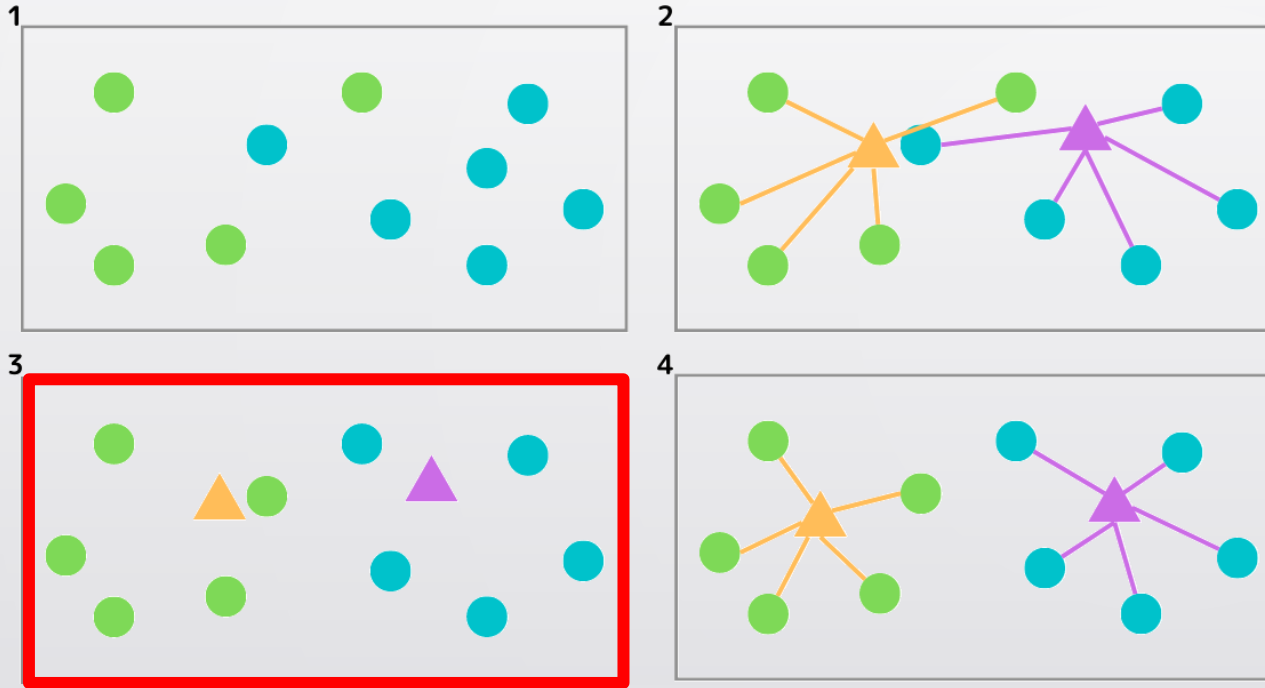


2. 各クラスターの中心を計算する
Compute the center of each cluster

3. 各クラスター中心とデータとの距離を計算する

Compute the distance of data from center of each cluster

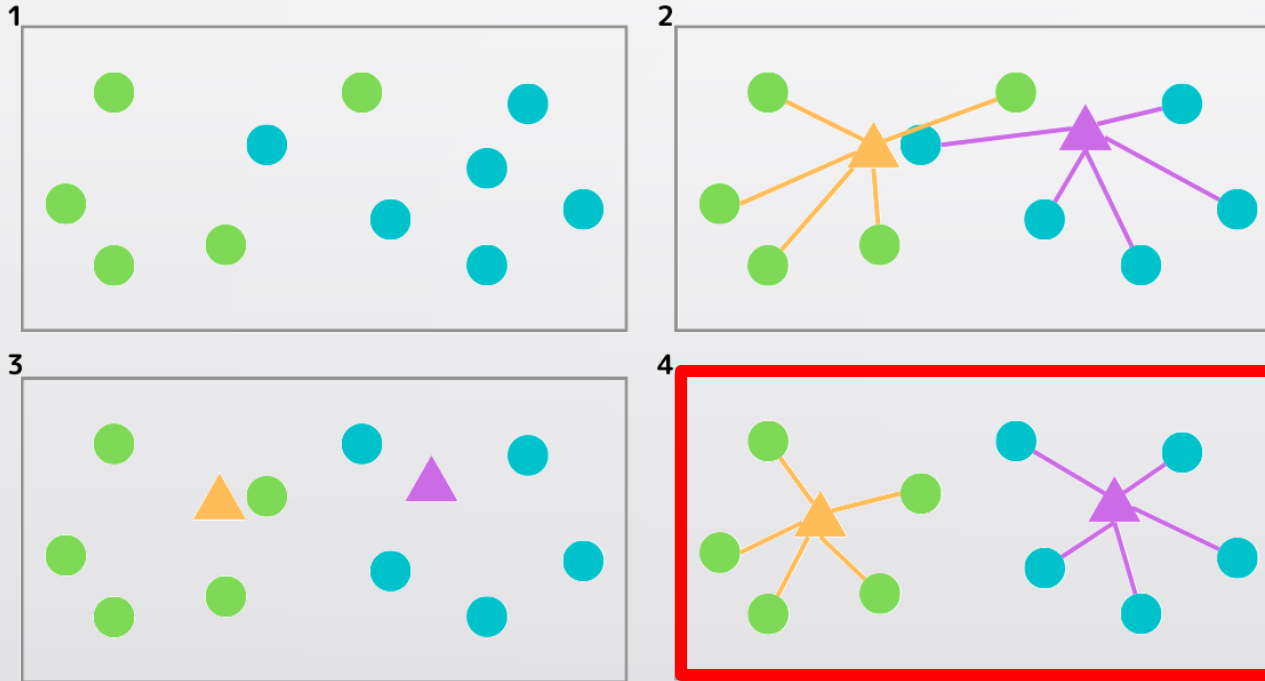
K平均クラスタリング k-means clustering



3. データと最も中心からの距離が近い
クラスターに割り当てる

Assign data point to the cluster with
smallest distance

K平均クラスタリング k-means clustering



4. データと最も中心からの距離が近い
クラスターに割り当てる

Assign data point to the cluster with
smallest distance

K平均クラスタリング k-means clustering

1. μ_k を固定し以下の方法で $q_{i,k}$ を決定する

Fix μ_k and determine $q_{i,k}$ following the rule below

$$q_{i,k} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|\mathbf{x}_i - \mu_j\|^2 \\ 0 & \end{cases}$$

\mathbf{x}_i を重心 μ_j との距離が一番近いクラスに割り当てる

Assign \mathbf{x}_i to the cluster whose centroid μ_j is closest to \mathbf{x}_i

2. μ_k を最適化する

Optimize μ_k

$$\mu_k = \frac{\sum_{i=1}^N q_{i,k} \mathbf{x}_i}{\sum_{i=1}^N q_{i,k}}$$

K平均クラスタリング K-means clustering

1-2.の手続きを収束するまで繰り返す Repeat Step 1-4 until the result converges

クラスター内の誤差平方和が閾値以下になることが収束条件である

Convergence Criterion is usually that squared-sum within cluster SSE_k becomes smaller than threshold

$$SSE_k = \sum \|x_i - \mu_k\|^2 \quad \sum SSE_k \leq Threshold$$

距離は、通常、ユークリッド距離を計算する

Usually, Euclidian distance is computed as the index of distance between data point and cluster center

クラスタ数決定法: エルボー法

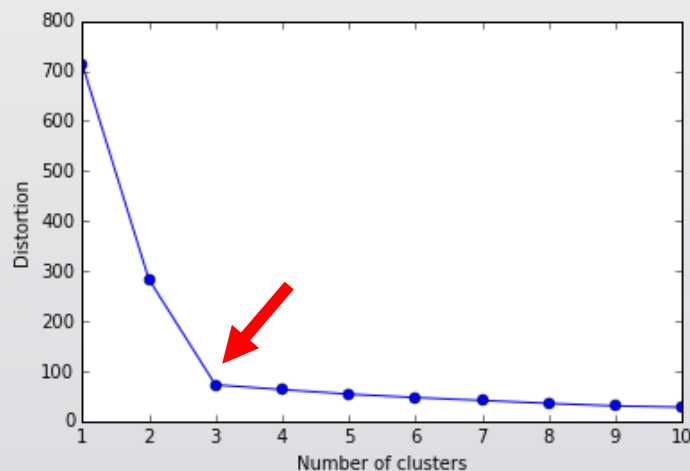
How to determine the number of Clusters: Elbow Method

1. 様々なクラスタ数でクラスタリングを行いSSEを計算する

Compute SSE after clustering with varying number of clusters

2. SSEをプロットし、SSEの減少が平坦になるクラスタ数を見つける

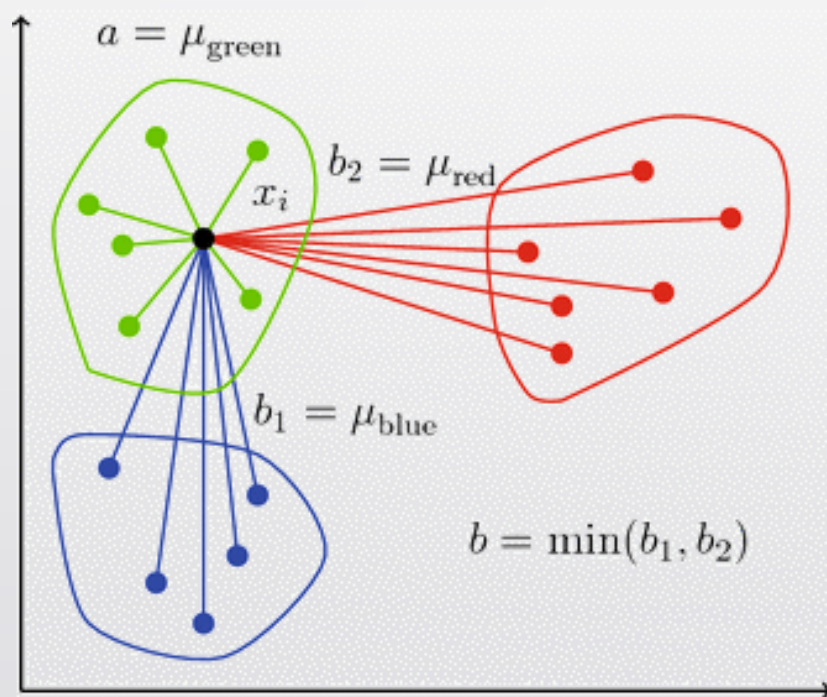
Plot SSEs and find the cluster number at which decrease of SSE reaches plateau



<https://qiita.com/deaikei/items/11a10fde5bb47a2cf2c2>

クラスタ数決定法: シルエット係数

How to determine the number of Clusters: Silhouette Coefficient



Usha Narra et al, 2016

a: クラスタ内の他のデータとの距離の平均
a: Mean distance from other data points within cluster

b: 最も近いクラスタのデータとの距離の平均
a: Mean distance from data points in nearest cluster

$$\text{Silhouette Coefficient} = \frac{b - a}{\max(b, a)}$$

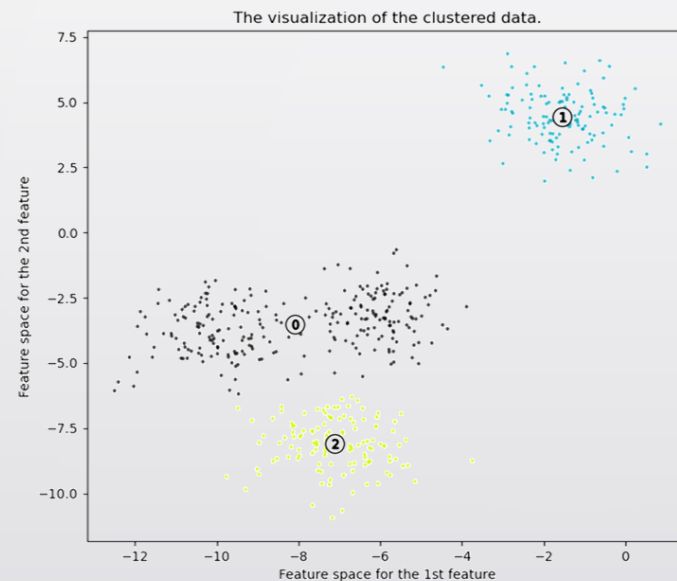
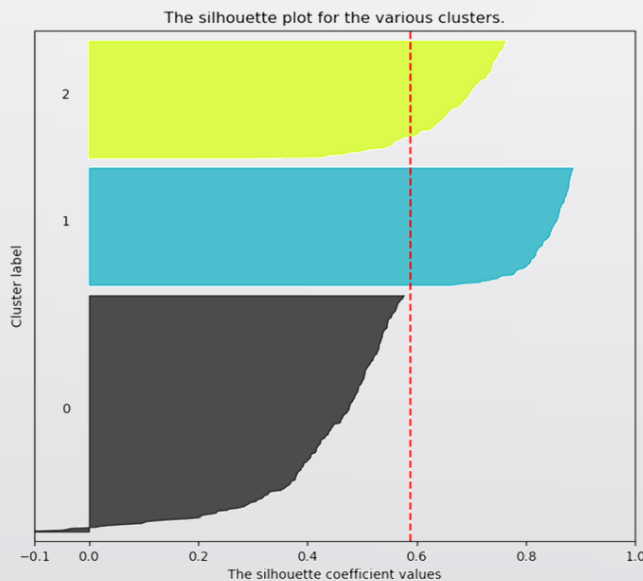
$[-1, 1]$ の範囲で変動する

Ranges within $[-1, 1]$

クラスタ数決定法: シルエット係数

How to determine the number of Clusters: Silhouette Coefficient

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 3$



全てのデータのシルエット係数をプロットしている
Silhouette coefficient of every data point is plotted

クラスタ数が不適切な時

When number of clusters is not appropriate

シルエット係数が小さなクラスタがある

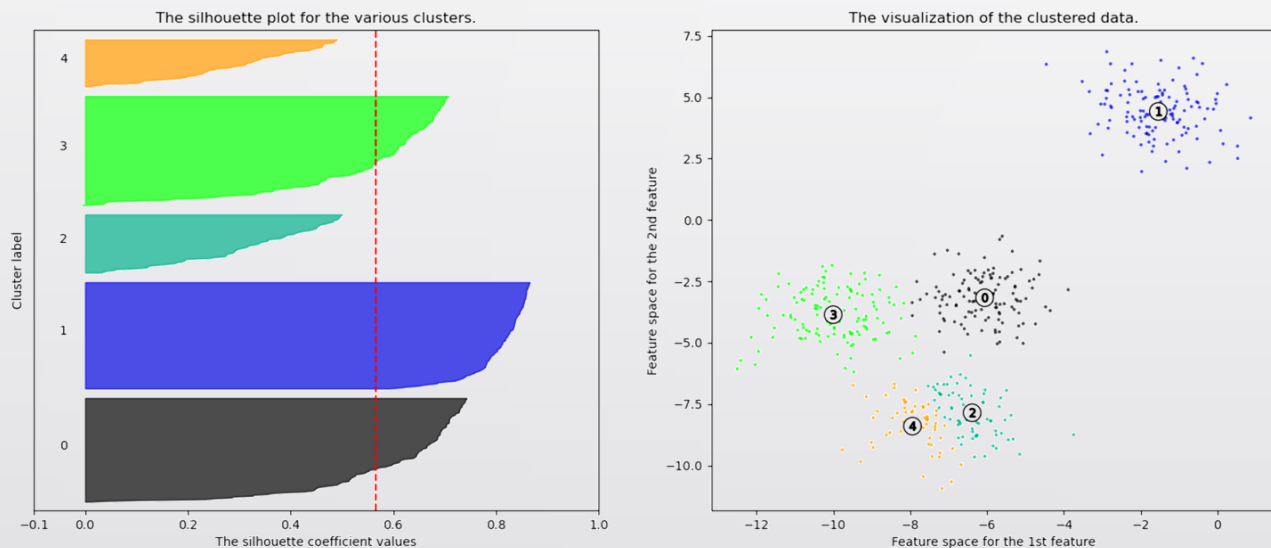
There are clusters with small silhouette coefficient

https://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

クラスタ数決定法: シルエット係数

How to determine the number of Clusters: Silhouette Coefficient

Silhouette analysis for KMeans clustering on sample data with `n_clusters = 5`



全てのデータのシルエット係数をプロットしている
Silhouette coefficient of every data point is plotted

クラスタ数が不適切な時

When number of clusters is not appropriate

シルエット係数が小さなクラスターがある

There are clusters with small silhouette coefficient

クラスタの大きさが不均一

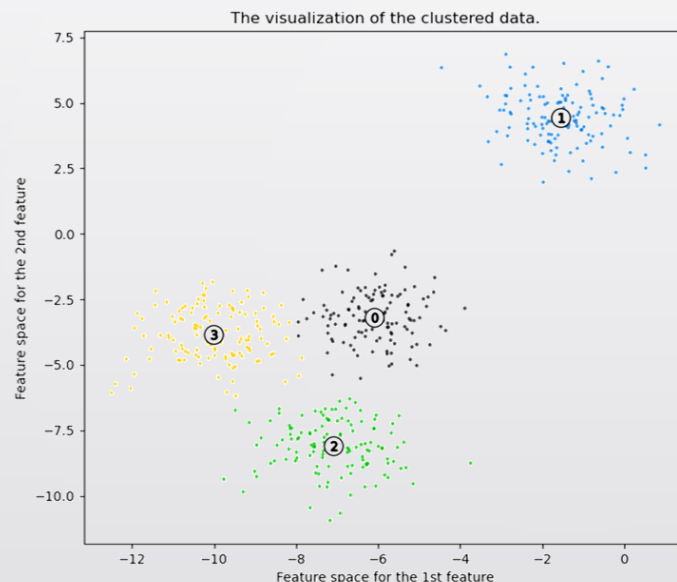
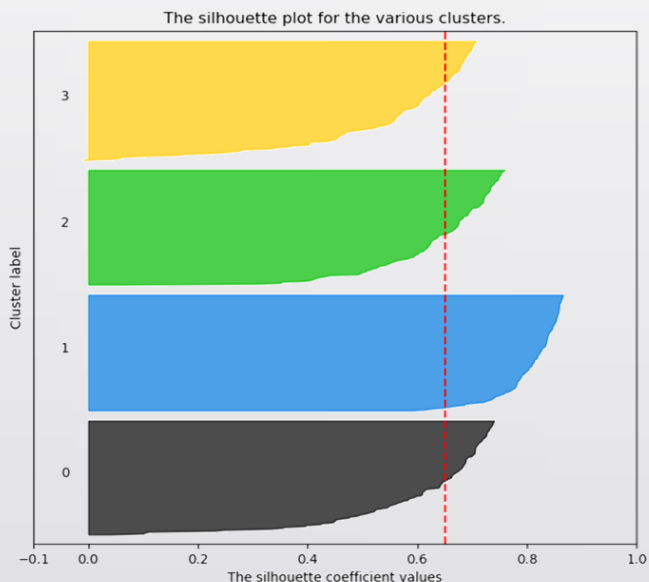
Size of cluster is inhomogenous

https://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

クラスタ数決定法: シルエット係数

How to determine the number of Clusters: Silhouette Coefficient

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 4$



クラスタ数が適切な時

When number of clusters is not appropriate

全てのクラスタのシルエット係数が大きい

Silhouette coefficient of every cluster is large enough

クラスタの大きさが均一

Size of cluster is homogenous

全てのデータのシルエット係数をプロットしている
Silhouette coefficient of every data point is plotted

https://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

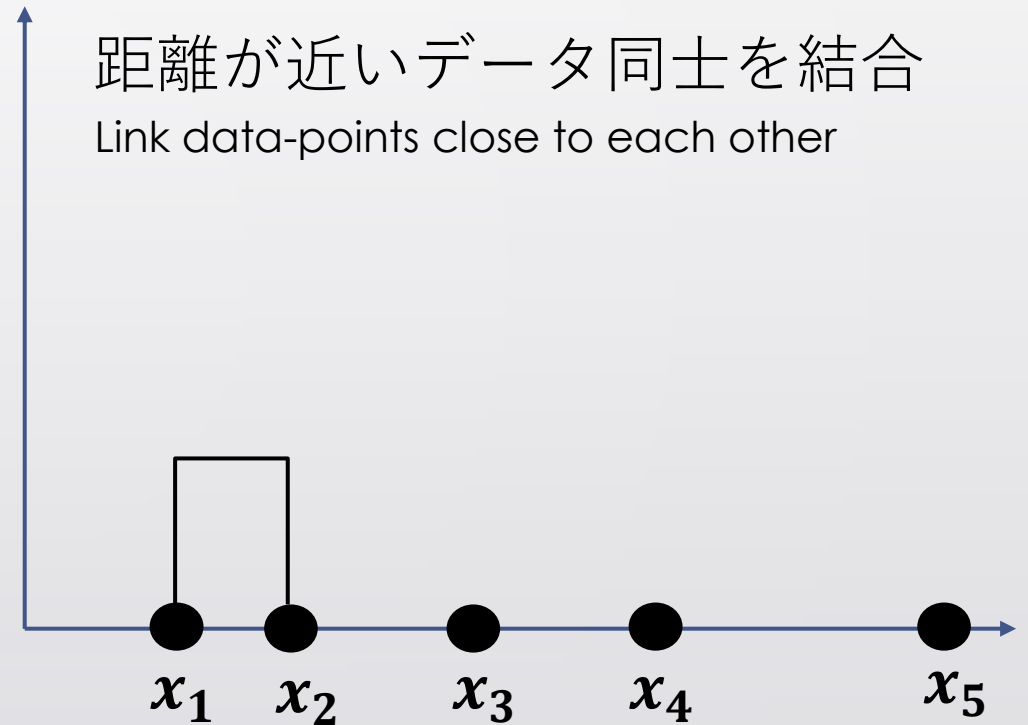
凝集性階層的クラスタリング

Agglomerative Hierarchical Clustering

	x_1	x_2	x_3	x_4	x_5
x_1	0				
x_2	1	0			
x_3	4	3	0		
x_4	6	5	2	0	
x_5	10	9	6	4	0

x_1	x_2	x_3	x_4	x_5
1	2	5	7	11

距離が近いデータ同士を結合
Link data-points close to each other

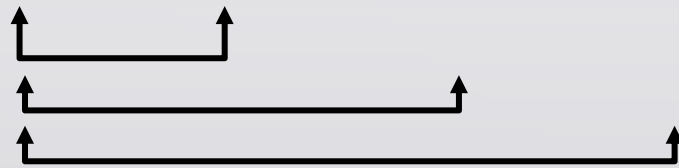


凝集性階層的クラスタリング

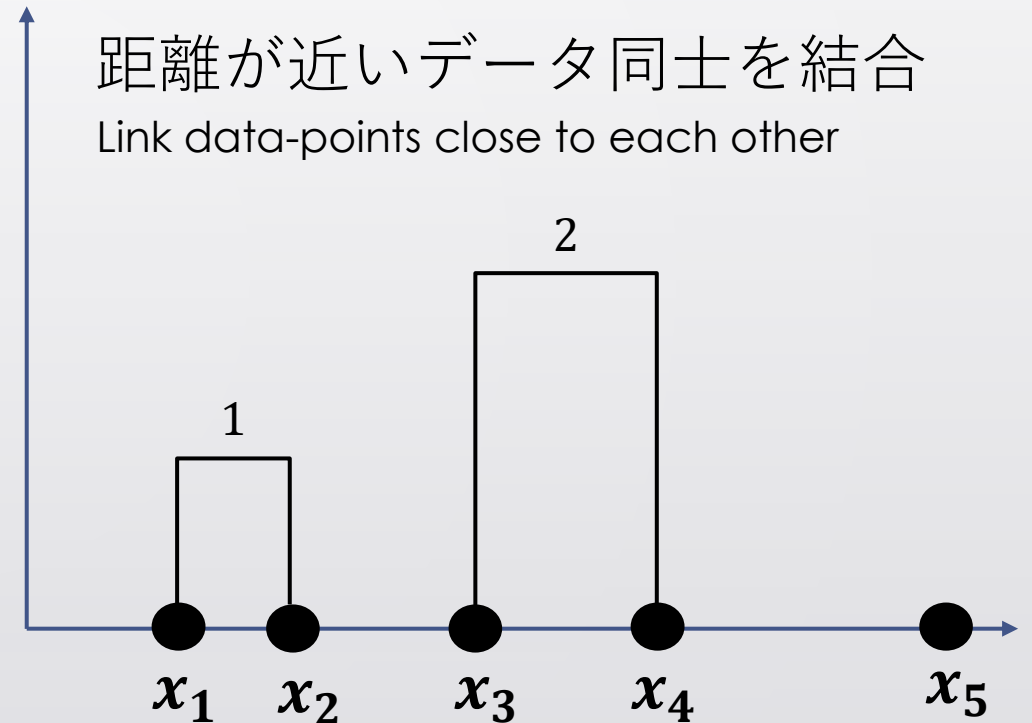
Agglomerative Hierarchical Clustering

	$\{x_1, x_2\}$	x_3	x_4	x_5
$\{x_1, x_2\}$	0			
x_3	3	0		
x_4	5	2	0	
x_5	9	6	4	0

x_1	x_2	x_3	x_4	x_5
1	2	5	7	11



距離が近いデータ同士を結合
Link data-points close to each other

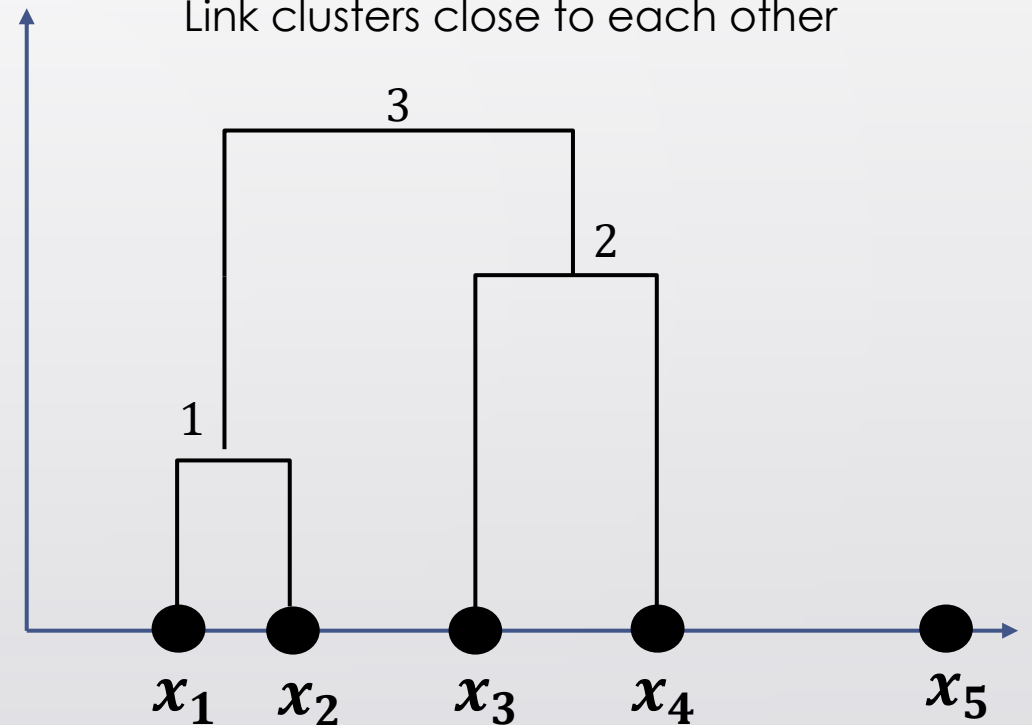
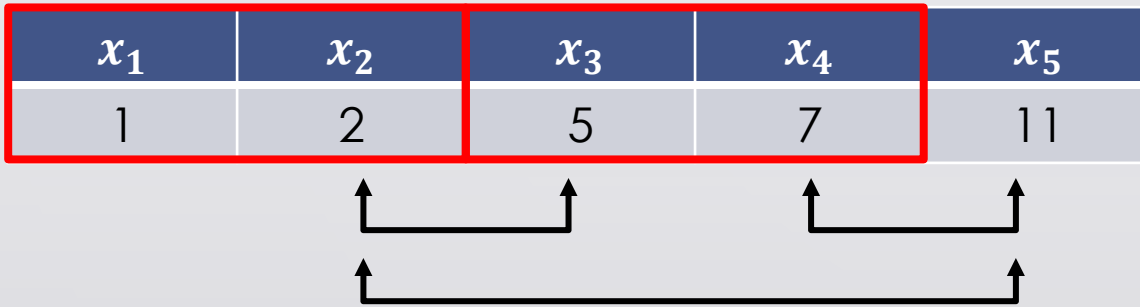


凝集性階層的クラスタリング

Agglomerative Hierarchical Clustering

距離が近いクラスター同士を結合
Link clusters close to each other

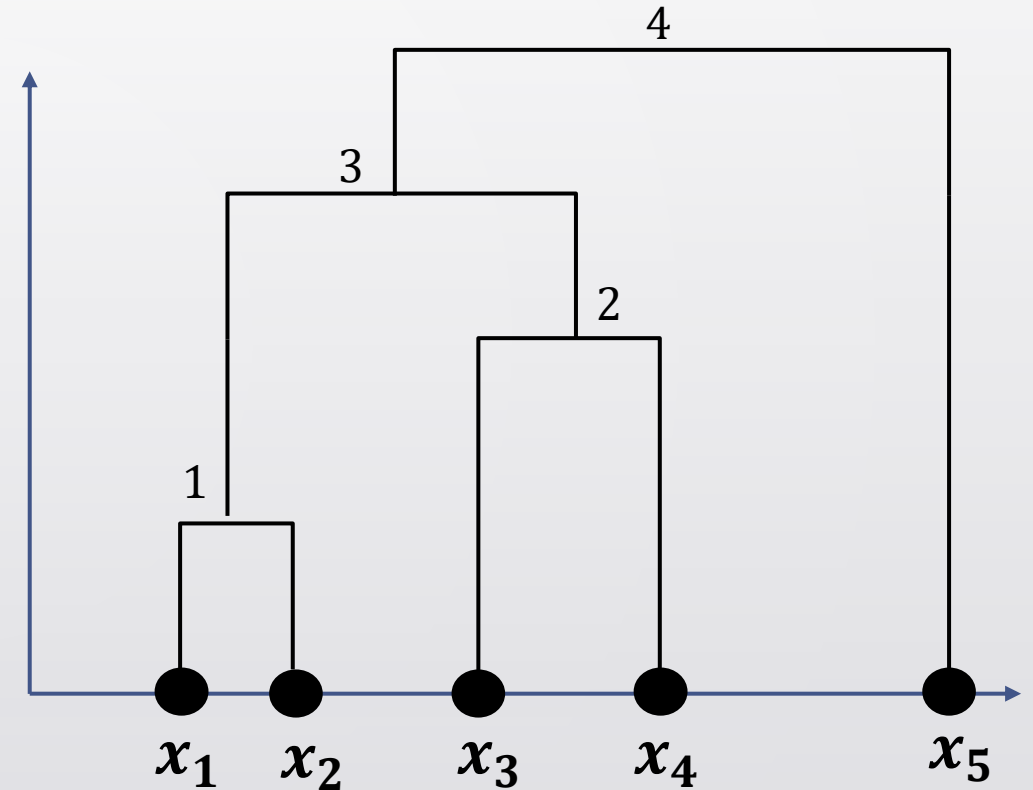
	$\{x_1, x_2\}$	$\{x_3, x_4\}$	x_5
$\{x_1, x_2\}$	0		
$\{x_3, x_4\}$	3	0	
x_5	9	4	0



凝集性階層的クラスタリング Agglomerative Hierarchical Clustering

	$\{x_1, x_2, x_3, x_4\}$	x_5
$\{x_1, x_2, x_3, x_4\}$	0	
x_5	4	0

x_1	x_2	x_3	x_4	x_5
1	2	5	7	11

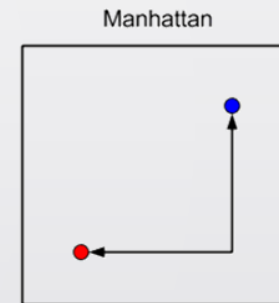


データ間の距離 Distance between Data-Points

ミンコフスキ距離 Minkowski Distance $Minkowski\ Distance = \left(\sum |x_{i,k} - x_{j,k}|^p \right)^{\frac{1}{p}}$

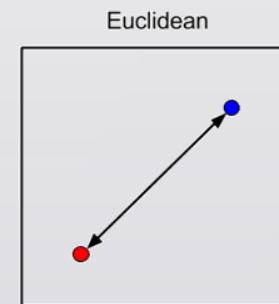
$p = 1, q = 1$ の時 In case of $p = 1, q = 1$

マンハッタン距離 $Manhattan\ Distance = \sum |x_{i,k} - x_{j,k}|$



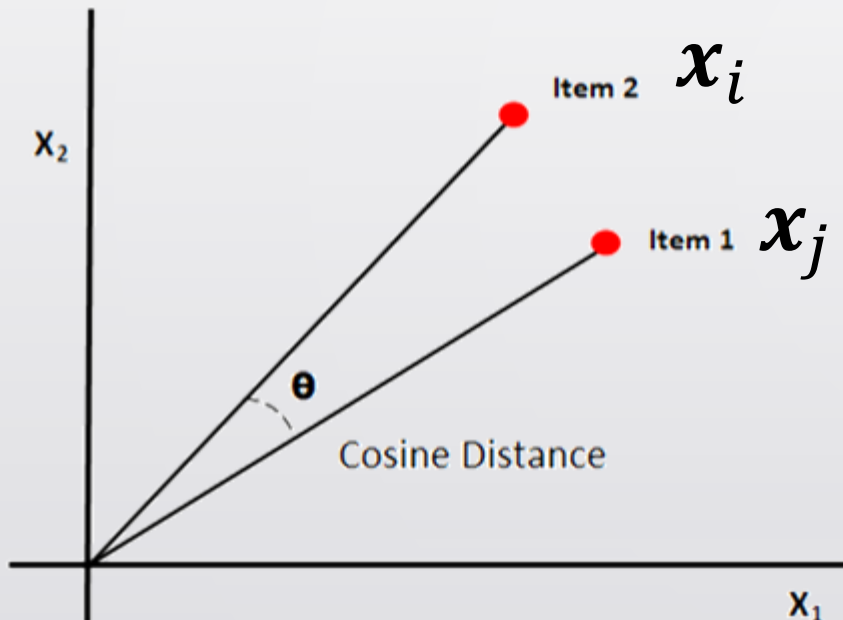
$p = 2, q = 2$ の時 In case of $p = 2, q = 2$

ユークリッド距離 $Euclidean\ Distance = \sqrt{\sum (x_{i,k} - x_{j,k})^2}$



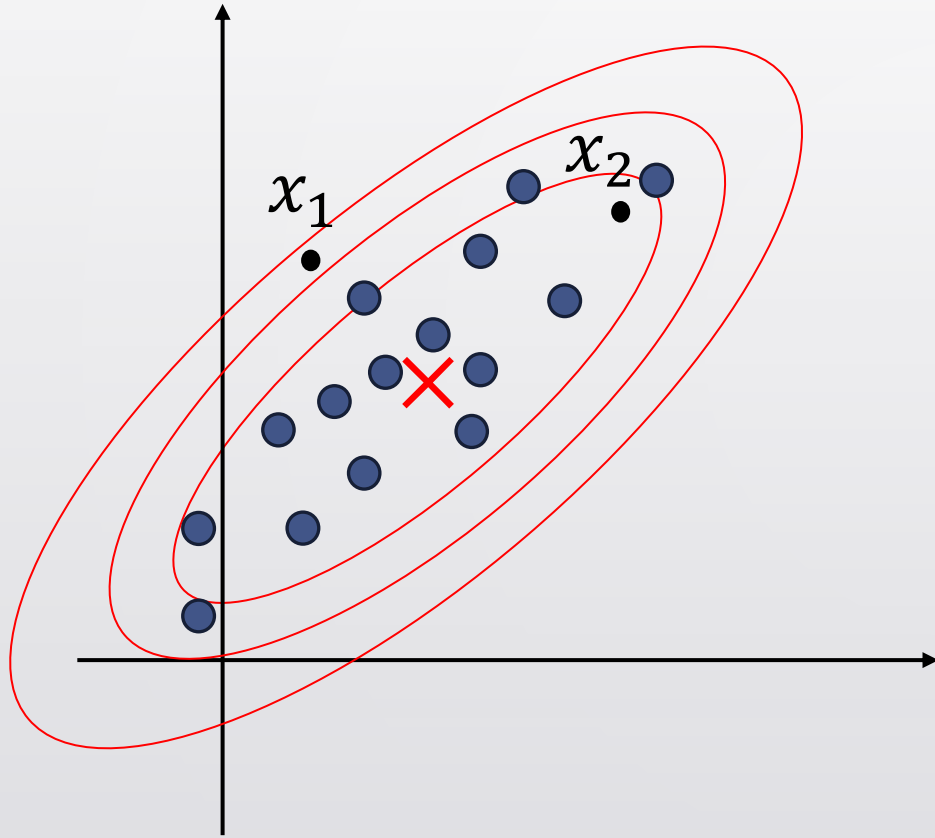
データ間の距離 Distance between Data-Points

コサイン類似度 Cosine Similarity



$$D(x_i, x_j) = \frac{\text{dot}(x_i, x_j)}{\|x_i\| \|x_j\|} = \frac{\sum x_{i,k} \times x_{j,k}}{\sum x_{i,k}^2 \times \sum x_{j,k}^2}$$

マハラノビス距離 Mahalanobis Distance



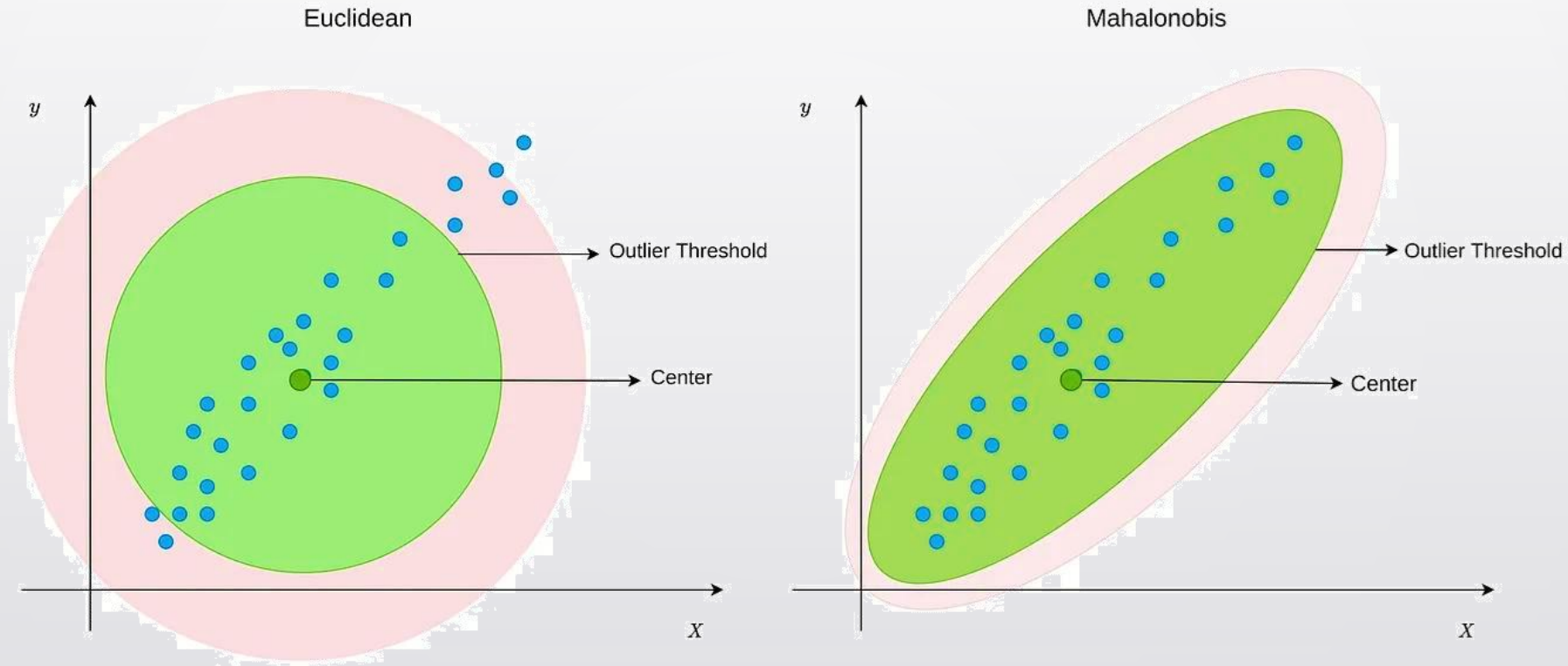
$$D_M(\mathbf{x}_i) = \sqrt{(\mathbf{x}_i - \bar{\mathbf{x}})\Sigma^{-1}(\mathbf{x}_i - \bar{\mathbf{x}})^T}$$

Σ : 分散共分散行列
Variance-Covariance Matrix

分布の形状を考慮した分布の中心からの距離の指標

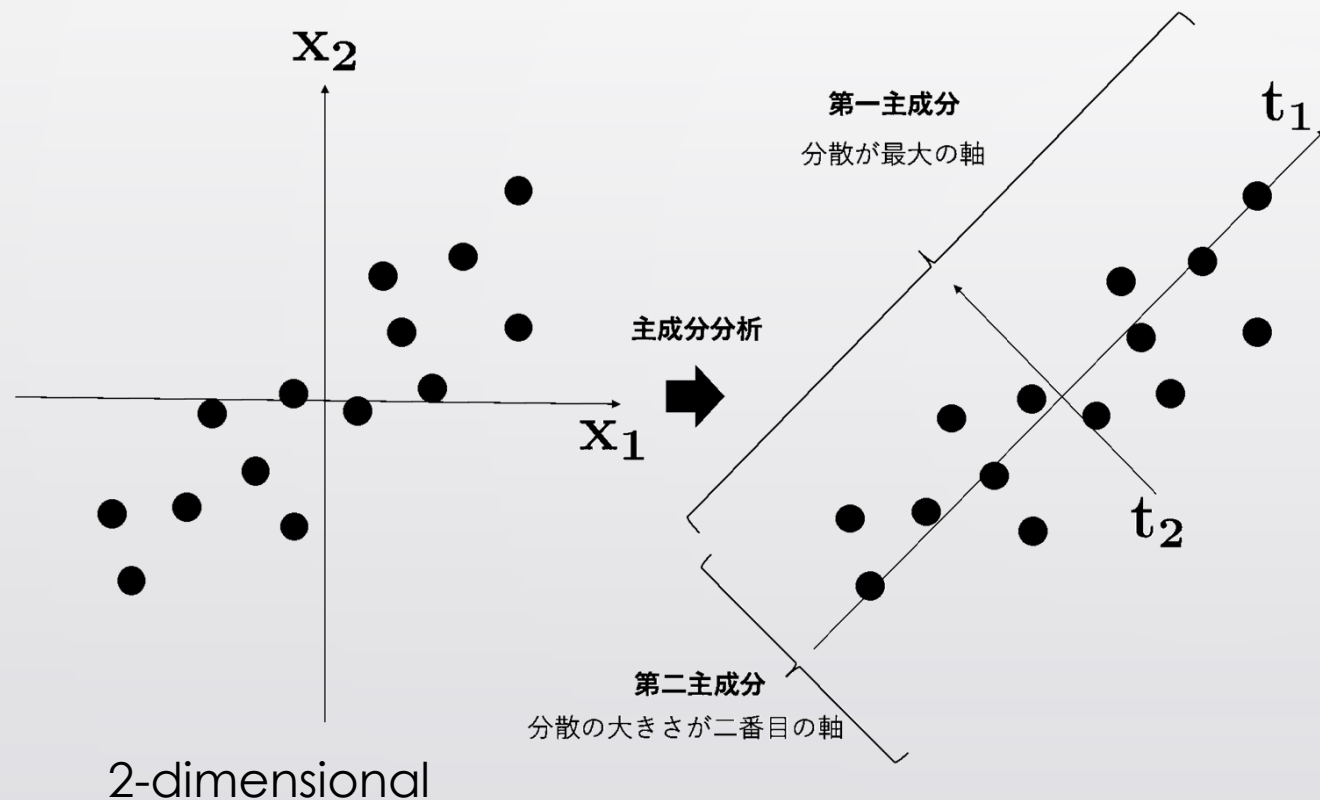
Measure of distance from distribution center adjusted by the shape of data distribution

マハラノビス距離 Mahalanobis Distance



<https://bob3.hatenablog.com/entry/2023/04/22/113540>

主成分 Principal Components



第1主成分軸は、データの分散が最大化される方向を向いている

The first PC axis is oriented in the direction along which variance of projected data is maximized

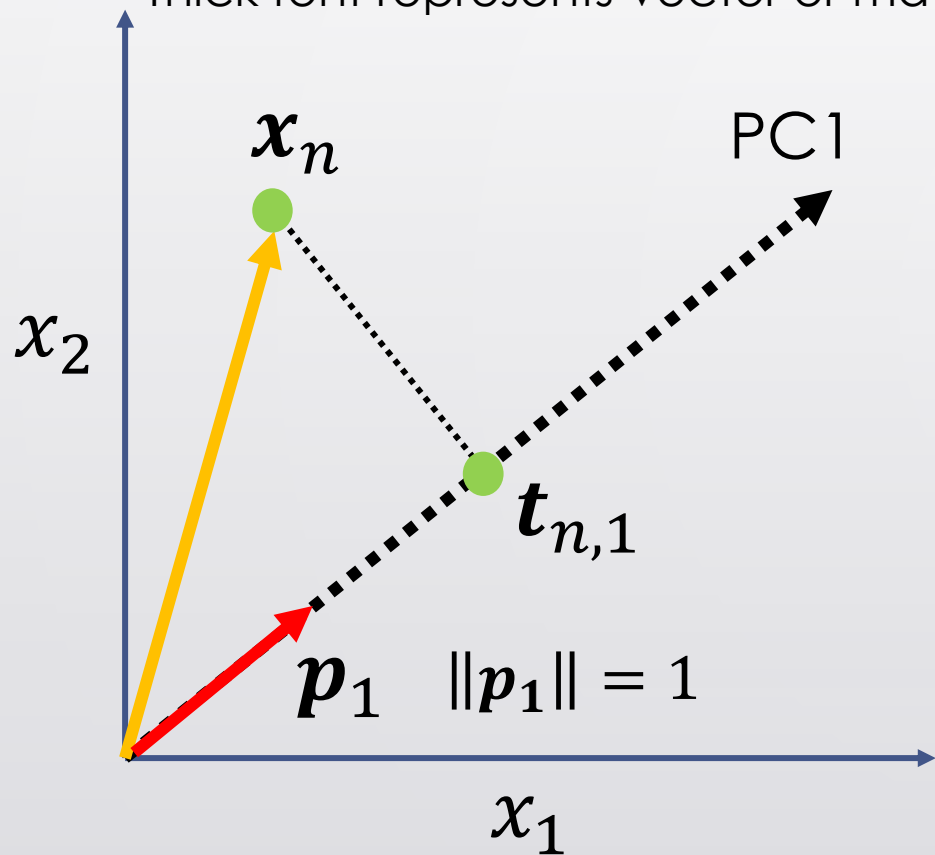
第 j 主成分軸は、データの分散が j 番目の大きさになる方向を向いている

The j -th PC axis is oriented in the direction along which projected data has j -th largest variance

第1主成分の計算 Computation of PC1

太字はベクトルか行列

Thick font represents vector or matrix



変数を中心化しておく Center the variables

観測データ x_n の第1主成分軸方向への射影 $t_{n,1}$ を計算する

Compute the projection $t_{n,1}$ of the observed data x_n onto the first PC axis

$t_{n,1}$ は x_n と p_1 の内積

$t_{n,1}$ is dot(inner) product of x_n and p_1

第1主成分の計算 Computation of PC1

$\mathbf{V}\mathbf{p}_1 = \lambda\mathbf{p}_1$ \mathbf{p}_1 は \mathbf{V} の固有ベクトルである \mathbf{p}_1 is eigenvector of \mathbf{V}

$\mathbf{V} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$ \mathbf{t}_1 の分散は λ に一致する Variance of \mathbf{t}_1 equals to λ

\mathbf{V} とは何か? What is \mathbf{V} ?

$$\mathbf{V} = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad \mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,M} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \cdots & x_{N,M} \end{bmatrix}$$

分散共分散行列 Variance-Covariance Matrix

\mathbf{V} は \mathbf{X} の分散共分散行列である \mathbf{V} is variance-covariance matrix of \mathbf{X}

$$\mathbf{V} = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \begin{bmatrix} x_{1,1} & x_{2,1} & \cdots & x_{N,1} \\ x_{1,2} & x_{2,2} & \cdots & x_{N,2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1,M} & x_{2,M} & \cdots & x_{N,M} \end{bmatrix} \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,M} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \cdots & x_{N,M} \end{bmatrix}$$
$$= \begin{bmatrix} \sigma_{1,1}^2 & \sigma_{1,2}^2 & \cdots & \sigma_{1,M}^2 \\ \sigma_{2,1}^2 & & \ddots & \vdots \\ \vdots & & & \sigma_{M-1,M}^2 \\ \sigma_{M,1}^2 & \sigma_{M,2}^2 & \cdots & \sigma_{M,M}^2 \end{bmatrix} \quad \sigma_{i,j}^2 = \frac{1}{N} \sum_{k=1}^N x_{k,i} x_{k,j}$$

分散共分散行列の対角化

Diagonalization of Variance-Covariance Matrix

$$\mathbf{P} = (\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M)$$

$$\mathbf{V}\mathbf{p}_i = \lambda_i\mathbf{p}_i \quad \|\mathbf{p}_i\| = 1$$

対称行列の異なる固有値に対する固有ベクトルは直交するので

Since eigenvectors of symmetric matrix corresponding to different eigen values are orthogonal

$$\mathbf{p}_i\mathbf{p}_j^T = \begin{cases} 1 & (\lambda_i = \lambda_j) \\ 0 & (\lambda_i \neq \lambda_j) \end{cases}$$

分散共分散行列の対角化

Diagonalization of Variance-Covariance Matrix

$$\mathbf{VP} = \mathbf{V}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M) = (\lambda_1 \mathbf{p}_1, \lambda_2 \mathbf{p}_2, \dots, \lambda_M \mathbf{p}_M)$$

$$\mathbf{p}_i^T \mathbf{VP} = (\lambda_1 \mathbf{p}_i^T \mathbf{p}_1, \lambda_2 \mathbf{p}_i^T \mathbf{p}_2, \dots, \lambda_M \mathbf{p}_i^T \mathbf{p}_M) = (0, 0, \dots, \lambda_i, \dots, 0)$$

$$\mathbf{P}^T \mathbf{VP} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \lambda_M \end{bmatrix}$$

マハラノビス距離の意味 Meaning of Mahalanobis Distance

$$\mathbf{x}_1 = (x_{11}, x_{12}) \quad \mathbf{x}_2 = (x_{21}, x_{22})$$

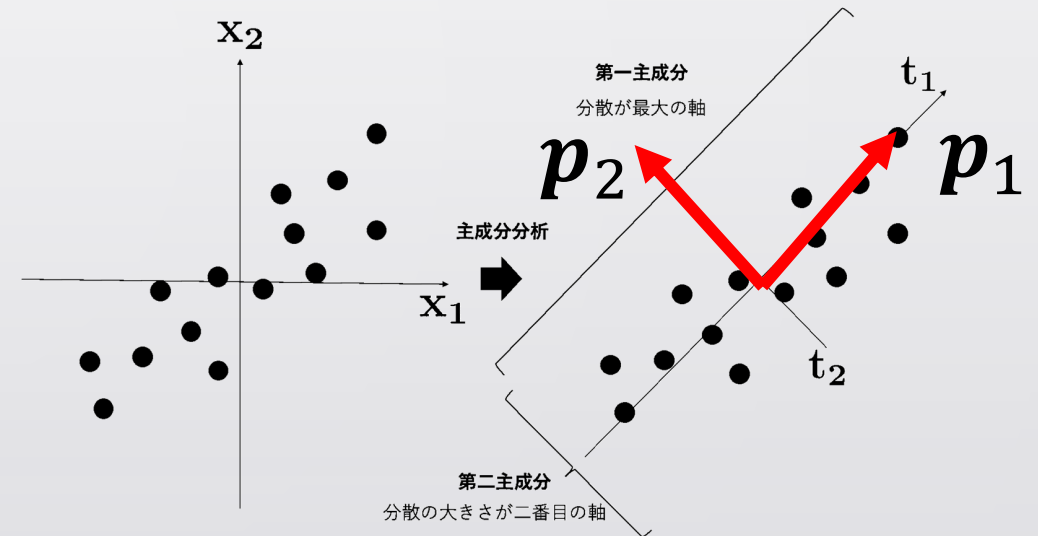
↓ \mathbf{x}_i の \mathbf{p}_k への射影を計算 Project \mathbf{x}_i onto \mathbf{p}_k

$$\mathbf{u}_1 = (u_{11}, u_{12}) = \mathbf{x}_1(\mathbf{p}_1, \mathbf{p}_2)$$

$$\mathbf{u}_2 = (u_{21}, u_{22}) = \mathbf{x}_2(\mathbf{p}_1, \mathbf{p}_2)$$

$(\mathbf{p}_1, \mathbf{p}_2)$ を基底とする座標系では \mathbf{x}_1 は \mathbf{u}_1 と表現される

\mathbf{x}_1 is expressed as \mathbf{u}_1 in a coordinate system defined by basis vectors of $(\mathbf{p}_1, \mathbf{p}_2)$



マハラノビス距離の意味 Meaning of Mahalanobis Distance

$$\mathbf{u}_1 = (u_{11}, u_{12}) = \mathbf{x}_1(\mathbf{p}_1, \mathbf{p}_2) \quad \mathbf{u}_2 = (u_{21}, u_{22}) = \mathbf{x}_2(\mathbf{p}_1, \mathbf{p}_2)$$

$$\begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1 \mathbf{p}_1 & \mathbf{x}_1 \mathbf{p}_2 \\ \mathbf{x}_2 \mathbf{p}_1 & \mathbf{x}_2 \mathbf{p}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} (\mathbf{p}_1, \mathbf{p}_2)$$

$$U = \begin{bmatrix} u_{1,1} & u_{1,2} & \cdots & u_{1,M} \\ u_{2,1} & u_{2,2} & \cdots & u_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ u_{N,1} & u_{N,2} & \cdots & u_{N,M} \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,M} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \cdots & x_{N,M} \end{bmatrix} (\mathbf{p}_1, \mathbf{p}_2 \cdots \mathbf{p}_M) = XP$$

マハラノビス距離の意味 Meaning of Mahalanobis Distance

$$U = \begin{bmatrix} u_{1,1} & u_{1,2} & \dots & u_{1,M} \\ u_{2,1} & u_{2,2} & \dots & u_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ u_{N,1} & u_{N,2} & \dots & u_{N,M} \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,M} \\ x_{2,1} & x_{2,2} & \dots & x_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,M} \end{bmatrix} (\mathbf{p}_1, \mathbf{p}_2 \dots \mathbf{p}_M) = X\mathbf{P}$$

U の分散共分散行列は Variance-covariance matrix Σ of U is

$$\Sigma = \frac{1}{N} \mathbf{U}^T \mathbf{U} = \frac{1}{N} (\mathbf{X}\mathbf{P})^T \mathbf{X}\mathbf{P} = \frac{1}{N} \mathbf{P}^T \mathbf{X}^T \mathbf{X} \mathbf{P} = \mathbf{P}^T \mathbf{V} \mathbf{P} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \lambda_M \end{bmatrix}$$

マハラノビス距離の意味 Meaning of Mahalanobis Distance

固有ベクトルで線型変換した後、マハラノビス距離を計算する

Calculate Mahalanobis distance after linear transformation by eigen vectors

$$D_M(\mathbf{u}_i) = \sqrt{(\mathbf{u}_i - \bar{\mathbf{u}})\Sigma^{-1}(\mathbf{u}_i - \bar{\mathbf{u}})^T}$$

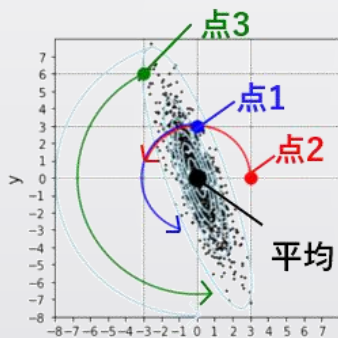
λ_i はデータの \mathbf{p}_i 方向の分散と一致

λ_i equals to variance of data along the direction of \mathbf{p}_i

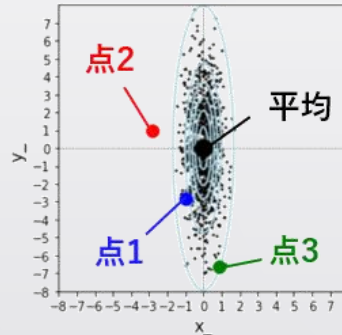
$$\Sigma^{-1} = \begin{bmatrix} \frac{1}{\lambda_1} & 0 & \dots & 0 \\ 0 & \frac{1}{\lambda_2} & & \vdots \\ & & \ddots & 0 \\ \vdots & & & \frac{1}{\lambda_M} \\ 0 & \dots & 0 & \frac{1}{\lambda_M} \end{bmatrix}$$

マハラノビス距離の意味 Meaning of Mahalanobis Distance

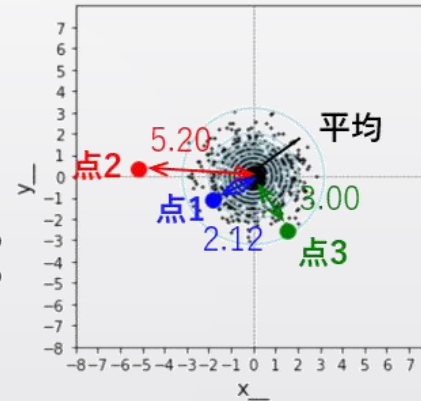
$$D_M(\mathbf{u}_i) = \sqrt{(\mathbf{u}_i - \bar{\mathbf{u}})\Sigma^{-1}(\mathbf{u}_i - \bar{\mathbf{u}})^T}$$



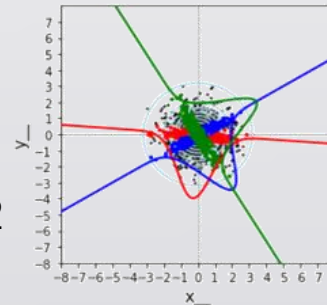
平均を中心に
約161度回転



横軸を0.55で割る
縦軸を2.59で割る



標準偏差1 標準偏差1
標準偏差1



すべての方向の分散を
一致させた後、ユーク
リッド距離を計算

Compute Euclidian
distance after equalizing
variance across all the
directions

固有ベクトルで線型変換
Linear transformation by eigen vectors

<https://qiita.com/yutera12/items/db425fafce2d87a25a1f>

Jaccard 係数 Jaccard Coefficient

集合の類似度の指標 Measure of similarity between two sets

ベクトル間の類似度の指標としても用いることができる

Can be used as a measure of similarity between vectors

$$x = \{x_1, x_2 \dots x_n\} \quad y = \{y_1, y_2 \dots y_n\}$$

$$Jaccard(x, y) = \frac{|x \cap y|}{|x \cup y|} \quad Jaccard(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}\mathbf{y}^T}{\mathbf{x}\mathbf{x}^T + \mathbf{y}\mathbf{y}^T - \mathbf{x}\mathbf{y}^T}$$

Dice 係数 Dice Coefficient

集合の類似度の指標 Measure of similarity between two sets

ベクトル間の類似度の指標としても用いることができる

Can be used as a measure of similarity between vectors

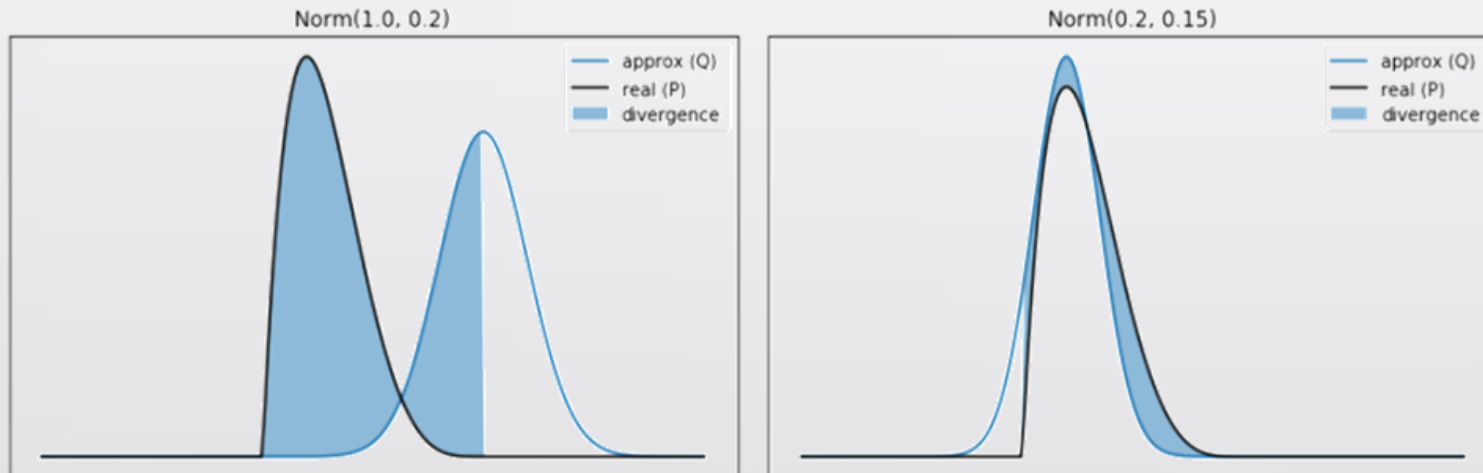
$$x = \{x_1, x_2 \dots x_n\} \quad y = \{y_1, y_2 \dots y_n\}$$

$$Dice(x, y) = \frac{2|x \cap y|}{|x| + |y|}$$

$$Dice(\mathbf{x}, \mathbf{y}) = \frac{2\mathbf{x}\mathbf{y}^T}{\|\mathbf{x}\| + \|\mathbf{y}\|}$$

KLダイバージェンス Kullback-Leibler Divergence

分布同士の類似度の評価指標 Measure of similarity between two distributions



$$KL(p(x)|q(x)) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right)$$

<https://jessicastringham.net/2018/12/27/KL-Divergence/>