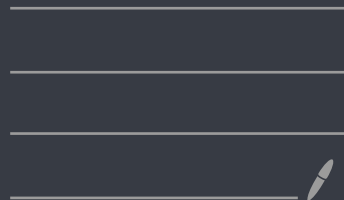


クラスター分析



クラスター分析 ...

似たものを集める

クラスター分析 ... 行方向にグルーピング

主成分、因子 ... 列方向にグルーピング

因子分析

id	use	design	price
1	3	2	3
2	3	5	1
3	4	1	4
⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮

$a_1 f_1 + a_2 f_2 + \epsilon$

クラスター分析

id	use	design	price
1	3	2	3
2	3	5	1
3	4	1	4
⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮

x_1
 x_3

都道府県

	人口	面積	人口密度
東京	;	.	.
大阪	,	,	,
北海道	,	,	,

主成分

... 変異度, 山地.
とかでデータ説明したい

ある量 L を最小にする

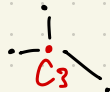
$$L(c_1, \dots, c_k) = \sum_{x \in \text{data}} \sum_{k=1}^k \|x - c_k\|_2^2$$

centroid

クラスターの重心
計算で出る

図にする

・ 3クラスター ... $k=3$



・ c_k とデータ点の距離をはかる.

ex. 東京 ... c_1 クラスター



$$d = \sqrt{(\text{東京人口} - c_1 \text{人口})^2 + (\text{東京面積} - c_1 \text{面積})^2 + (\text{東京人口密度} - c_1 \text{人口密度})^2}$$

※ ユークリッド距離

$$d = \|x - c_k\|_2 \quad \dots \quad \text{小さいほど, 距離が近い}$$

・ $L(c_1, \dots, c_k)$ は各データとそのクラスターの重心との距離を合計したもの

・ $L(c_1, \dots, c_k)$ が最小になるクラスター分けが良さそう.

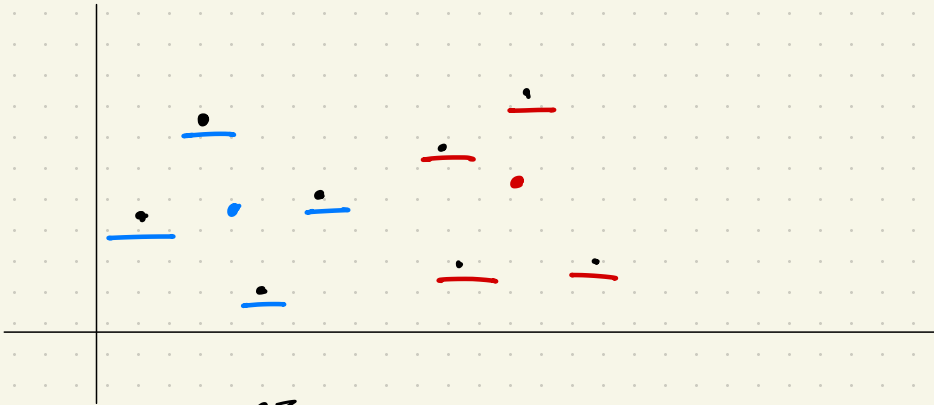
・ 厳密に計算するのは難しい, なるべく小さくなるように

$L(c_1 \dots c_k)$ を小さくしたい.

k-mean 法

- 1, クラスタ-所属をランダムに決める
- 2, 各クラスタ-の centroid が決まる \rightarrow 平均
- 3, 各データについて, 最も近い centroid を見つけ, クラスタ-を再設定
- 4, 2~3 を繰り返す
- 5, 再設定が最小になると終了. \rightarrow iter. max

$$\|x - c_k\|$$



k-mean 法の問題

\rightarrow 初期値がランダムなので, おかしなことになる!!

もうちょっとマシンに!!

k-mean ++ ... クラスタの重心って離れてるほうがいい!

1. 1つめの重心をデータ点の中からランダムに決める ... g_1

2. 2つめの重心は 1つめの重心から遠いほど選ばれやすくする.

・ 各データ点と g_1 の間: $\|x_i - g_1\|$

・ 選ばれた確率 $P(x) = \frac{(\|x_i - g_1\|)^2}{\sum (\|x_i - g_1\|)^2}$... g_2

3. 3つめの重心は、各データ点から g_1, g_2 の近い方を採用

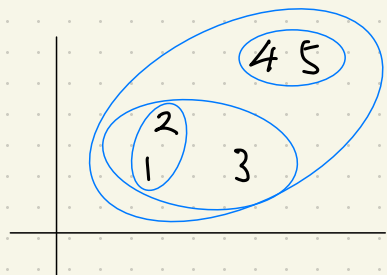
選ばれた確率を計算 ... g_1 から g_2 からも遠い点を選ばれた

4. 3を繰り返す.

真の最小値: L_{opt} とし、得られた L_{++} とした

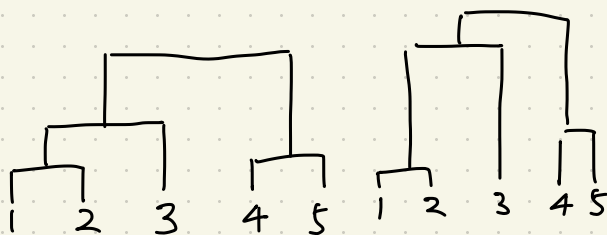
$$\underline{E(L_{++}) \leq 8(\log k + 2)L_{opt}}$$

階層クラスター分析

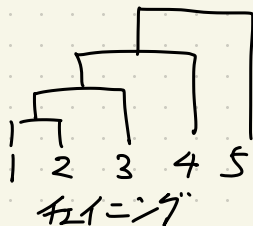


○で囲われているのがクラスター
あるクラスターが、大きなクラスターに含まれる。

↑
クラスター
間
の
距離



距離パターン

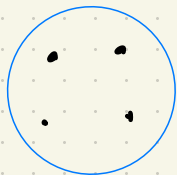
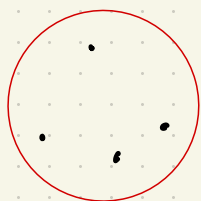


計算の順番

1. 各データ点を1つのクラスターと考える
2. 全てのクラスター間の距離を計算する。
3. 以下を繰り返す
 - a. 一番近いクラスターを併合
 - b. データ点が1つのクラスターにまとまったら計算終了。
 - c. 新しく作ったクラスターと既存のクラスターの距離を計算。

9. クラスター間の距離

- ・最短距離法
- ・最長距離法
- ・重心法
- ・Ward法

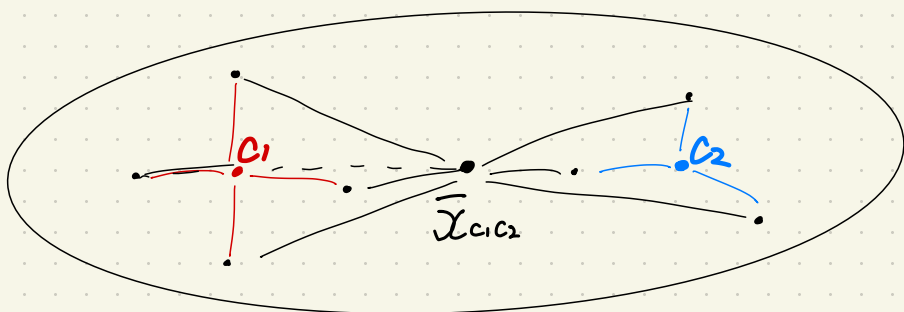


Ward 法

- ・ クラスター C_1, C_2 を合併すると仮定.
- ・ 合併後のクラスター内の重心からの相りの 2乗和を計算... ①
- ・ 合併前のクラスター内の重心からの相りの 2乗和を計算... ②
- ・ ① - ② が最小になるクラスターを合併する.

↓

クラスター内のバラツキを最小にするクラスターを合併する.



$$\sum \|x_i - \bar{x}_{C_1}\|^2$$

$$\sum \|x_i - \bar{x}_{C_2}\|^2$$

before

$$\sum \|x_i - \bar{x}_{C_1C_2}\|^2$$

after

$$\sum \|x_i - \bar{x}_{C_1C_2}\|^2 - \sum \|x_i - \bar{x}_{C_1}\|^2 - \sum \|x_i - \bar{x}_{C_2}\|^2$$

$$\Rightarrow \text{クラスター間の相り } d(C_1, C_2)$$

重心法の欠点
単調.

