

7. Cluster Analysis

樋口さぶろお

龍谷大学大学院理工学研究科数理情報学専攻

理論物理学特論 L09(2020-11-23 Mon)

最終更新: Time-stamp: "2020-11-21 Sat 15:10 JST hig"

今日の目標

- k -means clustering を行列で表現できる



クラスター分析とは KA MDA §7

$n \times p$ データ行列 X . KA MDA Table 6.2(p.87)

- 行数: n 学生数=職業数
- 列数: p 科目数=特性の個数
- 行列の要素: X_{ij} 学生 i の科目 j の点数, 職業 i の特性 j の値

学生・職業を点数・特性値が「近い」 k 個のグループ (cluster) に分けたい。

Membership Matrix G

$$G = (g_{ik}) \quad \boxed{\text{KA MDA (7.1) の 1 個前}} \quad \boxed{\text{KA MDA Table 7.1}}$$

- 行数: n 学生数=職業数
- 列数: K クラスターの個数
- 行列の要素: g_{ik} 属すれば 1, 属さなければ 0.

各行に成分 1 はちょうど 1 個だけ. 学生, 職業はちょうど 1 個のクラスターに属する.

$$\rightsquigarrow G \mathbf{1}_K = \mathbf{1}_n$$

Cluster Feature Matrix C

KA MDA §7.2

$$C = (c_{kj})$$

KA MDA (7.1) の 1 個前

KA MDA Table 7.2

- 行数: K クラスターの個数
- 列数: p 特性の個数
- 行列の要素: c_{ik} クラスターを代表する特性の値

Formulation KA MDA §7.2

KA MDA §7.2

$$X = GC + E \quad (7.3)$$

Minimize

$$f(G, C) = \|E\|^2 = \|X - GC\|^2 \quad (7.4)$$

K 次元散布図 KA MDA Fig.7.2

点: X の各行 (青四角) \tilde{x}'_i , C の各行 (赤丸) \tilde{c}'_k ,
 K 個の座標: K 個の各特性の値

Iterative Algorithm KA MDA 7.4

- explicit solution (陽解法) 解=既知の量の式 PCA
- implicit solution (陰解法) 解の満たす方程式がわかるだけ $f(\text{解})=0$
 - ▶ iterative solution (反復解法) だんだん真の解に収束していく (かも)

Algorithm (KMC= K -means Clustering) KA MDA p.98 下

- Step 1. $G_{[t]}, C_{[t]}$ を $t = 0$, 適当な可能解 $G_{[0]}, C_{[0]}$ にセット.
- Step 2. $f(G_{[t]}, C)$ を C について最小化, $C_{[t+1]}$ とおく.
- Step 3. $f(G, C_{[t+1]})$ を G について最小化, $G_{[t+1]}$ とおく.
- Step 4. 「収束したら」終了, そうでなければ Step 2. に戻る.

収束条件

収束条件の例「あまり変化しなくなった」

$$f(G_{[0]}, C_{[0]}) - f(G_{[1]}, C_{[1]}) \leq \epsilon$$

振動しないの? \rightsquigarrow

$$f(G_{[0]}, C_{[0]}) \geq f(G_{[0]}, C_{[1]}) \geq f(G_{[1]}, C_{[1]}) \geq f(G_{[1]}, C_{[2]}) \geq f(G_{[2]}, C_{[2]}) \geq \dots$$

真の最小値に到達する前に等号になっちゃわないの? \rightsquigarrow 収束先が最小値でない心配

global minimum 対 local minimum

Step 2. の具体的手続き KA MDA §7.5

$$f(G, C) = \|X - GC\|^2$$

を C に関して最小化したい.

KA MDA A.2.2 によれば

$$F(B) = \|Y - XB\|^2 \tag{A.2.11}$$

を最小化するには, X への射影 $P_{XY} = X(X'X)^{-1}X'Y$ を使って,

$$XB = P_X Y \tag{A.2.12}$$

となる, つまり $B = (X'X)^{-1}X'Y$ とせよ, とのことだった.

$Y = X, X = G, B = C$ として C を決めると,

$$C = (G'G)^{-1}G'X = D^{-1}G'X \tag{7.13}$$

Step 2. の直観的説明

$$C = (G'G)^{-1}G'X = D^{-1}G'X \quad (7.13)$$

D は対角行列で, 対角成分は, クラスタに所属する学生数.

$$\begin{aligned}
 c_{kj} = \bar{x}_{kj} &= \text{変数 } j \text{ の値 } x_{ij} \text{ の, クラスタ } k \text{ 内の平均値} \\
 &= \frac{1}{\sum_i g_{ik}} \sum_{i \in \text{cluster } k} x_{ij} \quad (7.15)
 \end{aligned}$$

Step 3. の具体的手続き KA MDA §7.6

$$f(G, C) = \|X - GC\|^2$$

を G に関して最小化したい。

行列にしたからきれいに書けるわけではない。 K 択が n 個あるだけだから大したことない。

いちばん小さくなる成分を 1 にしろ ... (7.19)

練習問題

KA MDA (7.5) の 10×2 のデータ行列 X を考える.

① $G_{[t]} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$ に対して, KA MDA §7.4 の Step 2 により, $C_{[t+1]}$ を求め

よう.

② $C_{[t]} = \begin{pmatrix} 7 & 3 \\ 4 & 1 \\ 6 & 7 \end{pmatrix}$ に対して, KA MDA §7.4 の Step 3 により, $G_{[t'+1]}$ を求めよう (3 行目まででいい).

Excel にコピーして計算することをおすすめ.