

## 10. 非階層的クラスター分析(k-means 法)

**概要:** 前節では因子の似ている者同士を評価し、クラスターを構成した。これを階層的クラスター分析と呼ぶ。ここでは、最初からクラスターの数指定し、データがどのクラスターに属するのかを評価する方法を示す。これをヒクラスター分析と呼ぶ。この方法においては、各グループの平均を評価し、それをグループの特徴とする。各データがどのグループに近いかの距離を評価し、一番近いグループに属させる。その後、そのグループごとの平均を評価し直し、各データとの距離を再評価し、どのグループに属するかを決定する。この操作を繰り返し、データの移動、つまり各グループの平均の変化が無くなるまで続ける。

**キーワード:** クラスター分析; 非階層的クラスター分析;  $k$ -means 法.

### 10.1. 序

前節では因子の似ている者同士を評価し、クラスターを構成した。ここでは、最初からクラスターの数指定し、データがどのクラスターに属するのかを評価する方法を示す。

### 10.2. k-means 法

ここでは  $k$ -means 法と呼ばれる手法を扱う。この方法では、まずクラスターの数設定する。このため、この方法のことを非階層的クラスター分析と呼ぶ。そして、平均値の初期値を設定する。ここでは、クラスターの数  $k$  とする。

我々は  $P$  種類の因子を考える。

初期値として、我々はそれぞれのクラスターのそれぞれの因子の平均値を以下のように設定する。

$$\mu_{1_{C_1}}, \mu_{2_{C_1}}, \dots, \mu_{p_{C_1}} \quad (1)$$

$$\mu_{1_{C_2}}, \mu_{2_{C_2}}, \dots, \mu_{p_{C_2}} \quad (2)$$

...

$$\mu_{1_{C_k}}, \mu_{2_{C_k}}, \dots, \mu_{p_{C_k}} \quad (3)$$

ここで、 $\mu_{i_{C_j}}$  は  $i$  因子のクラスター  $C_j$  における平均である。これらの平均値は仮のものであり、後に変更される。

上の平均値を使い、あるメンバーのデータ

$$\mathbf{x}_i = (x_{i1} \quad x_{i2} \quad \cdots \quad x_{ip}) \quad (4)$$

を評価する。

このデータと  $\mathbf{x}_i$  クラスタ  $j$  の距離  $D_{C_j}^2[\mathbf{x}_i]$  として以下を定義する。

$$D_{C_1}^2[\mathbf{x}_i] = (x_{i1} - \mu_{1_{-C_1}})^2 + (x_{i2} - \mu_{2_{-C_1}})^2 + \cdots + (x_{ip} - \mu_{p_{-C_1}})^2 \quad (5)$$

$$D_{C_2}^2[\mathbf{x}_i] = (x_{i1} - \mu_{1_{-C_2}})^2 + (x_{i2} - \mu_{2_{-C_2}})^2 + \cdots + (x_{ip} - \mu_{p_{-C_2}})^2 \quad (6)$$

...

$$D_{C_k}^2[\mathbf{x}_i] = (x_{i1} - \mu_{1_{-C_k}})^2 + (x_{i2} - \mu_{2_{-C_k}})^2 + \cdots + (x_{ip} - \mu_{p_{-C_k}})^2 \quad (7)$$

この中の最小のものを選択する。

$$\text{Min}[D_{C_1}^2[\mathbf{x}_i], D_{C_2}^2[\mathbf{x}_i], \dots, D_{C_k}^2[\mathbf{x}_i]] = D_{C_\zeta}^2[\mathbf{x}_i] \quad (8)$$

これにより、このデータは距離の二乗の最小値を与えるクラスタ  $\zeta$  に属すると判断する。この解析を全ての新しいデータに対して行う。すると、各データはどれかのクラスタに属することになる。

ここで、我々は各クラスタの平均値を再度評価する。

$$\mu_{1_{-C_1}} = \frac{\sum_{i=1}^{n_{C_1}} x_{i1}}{n_{C_1}}, \mu_{2_{-C_1}} = \frac{\sum_{i=1}^{n_{C_1}} x_{i2}}{n_{C_1}}, \dots, \mu_{p_{-C_1}} = \frac{\sum_{i=1}^{n_{C_1}} x_{ip}}{n_{C_1}} \quad \text{for } \mathbf{x}_i \in C_1 \quad (9)$$

$$\mu_{1_{-C_2}} = \frac{\sum_{i=1}^{n_{C_2}} x_{i1}}{n_{C_2}}, \mu_{2_{-C_2}} = \frac{\sum_{i=1}^{n_{C_2}} x_{i2}}{n_{C_2}}, \dots, \mu_{p_{-C_2}} = \frac{\sum_{i=1}^{n_{C_2}} x_{ip}}{n_{C_2}} \quad \text{for } \mathbf{x}_i \in C_2 \quad (10)$$

...

$$\mu_{1_{-C_k}} = \frac{\sum_{i=1}^{n_{C_k}} x_{i1}}{n_{C_k}}, \mu_{2_{-C_k}} = \frac{\sum_{i=1}^{n_{C_k}} x_{i2}}{n_{C_k}}, \dots, \mu_{p_{-C_k}} = \frac{\sum_{i=1}^{n_{C_k}} x_{ip}}{n_{C_k}} \quad \text{for } \mathbf{x}_i \in C_k \quad (11)$$

ここで、 $n_{C_\zeta}$  はクラスタ  $\zeta$  のデータ数である。

この新たな平均値を使い、我々は全てのデータを評価し直す。このプロセスをすべての平均値が変動しなくなるまで繰り返す。つまり、所属を変更するデータがなくなるまで繰り返す。

このプロセスを終えると、我々は最終的な  $k$  個のクラスタを得る。

### 10.3. k-means 法の初期設定

この方法には、初期のグループ数を決定することと、そのグループの平均の設定に任意性がある。

一つは、データをランダムに振り分ける方法である。これは、あるデータに対して乱数を発生させてどれかのグループに属させることをする。その後、クラスターを構成をする。

一つは、項目に対してある一定の値を  $k$  個設定し、それに対して各データを振り分ける方法である。この場合は、一定の値は最後までその値のままである。

もう一つは、初期の比較的少ないデータで階層的クラスター分析をおこない、クラスター数、平均値を設定する。その後は、これまで示してきた k-means 法を適用する。

### 10.4. まとめ

この章のまとめを行う。

$p$  個の因子があるとする。ここで、 $k$  個のクラスターがあると設定する。そして、クラスター  $\varsigma$  の仮の平均を以下のように設定する。

$$\mu_{1_{-C_\varsigma}}, \mu_{2_{-C_\varsigma}}, \dots, \mu_{p_{-C_\varsigma}}$$

ここで  $\varsigma = 1, 2, \dots, k$  である。

これらを使い、各データとクラスターの距離の二乗を以下のように評価する。

$$D_{C_\varsigma}^2[\mathbf{x}_i] = (x_{i1} - \mu_{1_{-C_\varsigma}})^2 + (x_{i2} - \mu_{2_{-C_\varsigma}})^2 + \dots + (x_{ip} - \mu_{p_{-C_\varsigma}})^2$$

この距離が最小のものをそのデータが属するクラスターとする。これを全てのデータで繰り返す。その後、平均値を以下のように再評価する。

$$\mu_{1_{-C_\varsigma}} = \frac{\sum_{i=1}^{n_{C_\varsigma}} x_{i1}}{n_{C_\varsigma}}, \mu_{2_{-C_\varsigma}} = \frac{\sum_{i=1}^{n_{C_\varsigma}} x_{i2}}{n_{C_\varsigma}}, \dots, \mu_{p_{-C_\varsigma}} = \frac{\sum_{i=1}^{n_{C_\varsigma}} x_{ip}}{n_{C_\varsigma}} \quad \text{for } \mathbf{x}_i \in C_\varsigma$$

ここで、 $n_{C_\varsigma}$  はクラスター  $\varsigma$  のデータ数である。

このプロセスを平均値が変動しなくなるまで繰り返す。最終的には、我々は  $k$  個のクラスターを得る。

この方法の初期設定に、3 種類の方法を示した。