

清水 信夫 データ科学研究系 助教

大量の関数データが存在するとき、これらの少数のグループによる分類や特徴的な関数データによる要約が必要となる場合がしばしば生じるが、そのような場合の解析手法として関数データ解析(Ramsay, 1982; Ramsay & Silverman, 1997, 2005)の一手法である関数クラスタリングがよく用いられている。

一方、 k -meansクラスタリングと同様の規準に基づいて与えられた確率変数に関し空間を k 個の領域に分割したときの重心として主要点(principal points)が定義されている(Flury, 1990)。関数データ解析においても関数データを確率的に扱う概念としてランダム関数(Ibragimov & Rozanov, 1978)が定義されており、それにおける主要点(関数主要点)(Tarpey & Kinatader, 2003)が提案されている。ここでは、正規ランダム関数に従う関数データ集合に対して k -means関数クラスタリングを適用した場合の各クラスターの重心と正規ランダム関数における関数主要点の関連、および k -means関数クラスタリングを行った場合における局所解の出現回数について示す(清水・水田(2008))。

主要点とは (Flury, 1990, 1993)

p 変量確率変数における主要点とは、確率分布を代表する p 次元空間の点の集合として定義される。すなわち、確率変数 X の分布関数 F について

$\zeta_j, y_j \in R^p$ ($1 \leq j \leq k$)、かつ A_j および D_j がそれぞれ

$$A_j = \{x \in R^p; (x - \zeta_j)'(x - \zeta_j) < (x - \zeta_i)'(x - \zeta_i), \forall i \neq j\}$$

$$D_j = \{x \in R^p; (x - y_j)'(x - y_j) < (x - y_i)'(x - y_i), \forall i \neq j\}$$

をみたすとき、任意の y_j に対して

$$\sum_{j=1}^k \Pr(X \in A_j) E\{(X - \zeta_j)'(X - \zeta_j) | X \in A_j\} \leq \sum_{j=1}^k \Pr(X \in D_j) E\{(X - y_j)'(X - y_j) | X \in D_j\}$$

が成り立つような $\{\zeta_1, \dots, \zeta_k\}$ は、 F に従う p 変量確率変数 X の主要点である。

なお、 A_j および D_j はそれぞれ ζ_j および y_j の Domain of Attraction である。

ランダム関数とは (Ibragimov & Rozanov, 1978)

確率空間 (Ω, B, P) において、 $L^2(T)$ を実数空間 T 上の2乗可積分かつ可測な全ての関数からなる完備ヒルベルト空間、 ν を Ω から L^2 への可測写像とすると、 ν の要素はランダム関数と呼ばれる。

なお、 $\xi \in L^2(T)$ を実数空間 T 上の連続なランダム関数、

$y_j \in L^2(T)$ ($1 \leq j \leq k$) を非ランダム関数とすると、 y_j の Domain of Attraction は

$$A_j = \{\xi \in L^2[T_L, T_U]; \|\xi - y_j\|^2 < \|\xi - y_i\|^2, \forall i \neq j\}$$

と定義される。

主要点と関数主要点との関係 (Tarpey & Kinatader, 2003)

関数データ解析においては、関数データを正規直交基底を用いて展開する手法が広く利用されている。この手法をランダム関数にも適用することで、

ランダム関数の関数主要点を、ランダム関数の正規直交基底の各係数を成分とする空間上における多変量確率変数の主要点に対応する関数として表すことができる。

関数主要点 (Tarpey & Kinatader, 2003)

$\zeta_j, y_j \in L^2(T)$ ($1 \leq j \leq k$) をそれぞれ非ランダム関数とし、 A_j および D_j をそれぞれ ζ_j および y_j の Domain of Attraction、 z および y をそれぞれ

$z = \sum_{j=1}^k \zeta_j I_{\{\xi \in A_j\}}$ および $y = \sum_{j=1}^k \zeta_j I_{\{\xi \in D_j\}}$ で表されるランダム関数とする。

ここで $\{\zeta_1, \dots, \zeta_k\}$ が ξ の関数主要点であるとは、任意の y_j に対して

$$E\|\xi - z\|^2 \leq E\|\xi - y\|^2$$

関数 k -meansクラスタリングアルゴリズム

$y(t) \in L^2[T_L, T_U]$: ランダム関数

$y_i(t) \in L^2[T_L, T_U]$: $y(t)$ の標本関数 (非ランダム関数) ($1 \leq i \leq n$)

$m_j(t)$ ($1 \leq j \leq k$): クラスター j における中心関数

n_j : クラスター j に含まれる標本関数の個数

$A_j = \{\xi \in L^2[T_L, T_U]; \|\xi - y_j\|^2 < \|\xi - y_i\|^2, \forall i \neq j\}$: $y_j(t)$ の Domain of Attraction

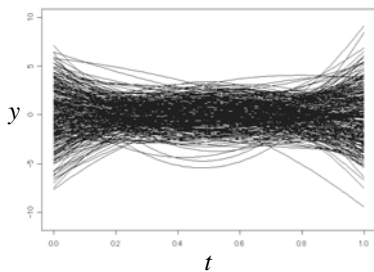
(1) $m_j(t)$ の初期値を N 組与える。

(2) 各 $y_i(t)$ について $\int_{T_L}^{T_U} (m_j(t) - y_i(t))^2 dt$ が最小となる j を求め、 $y_i(t)$ を A_j に振り分ける。

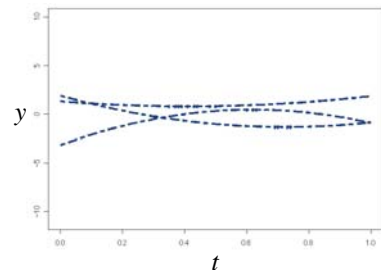
(3) $m_j(t) = \sum_{y_i(t) \in A_j} y_i(t) / n_j$ を計算する。

(4) (2) ~ (3) を繰り返し、 $m_j(t)$ が変化しなくなったときの値および A_j に含まれる関数データ集合を中心関数および各クラスターとする。

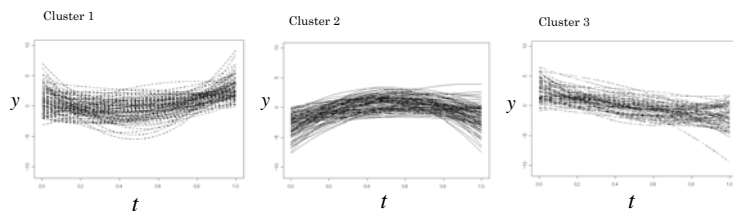
$y(t) = a_1 + a_2 t + a_3 t^2 = b_1 + \sqrt{12}(t-1/2)b_2 + \sqrt{12}(t^2-t+1/6)b_3$ ($0 < t < 1$) における関数クラスタリング ($(b_1, b_2, b_3)' \sim N_3((0,0,0)', \text{diag}(1,1,1))$) ($n=250, k=3$)



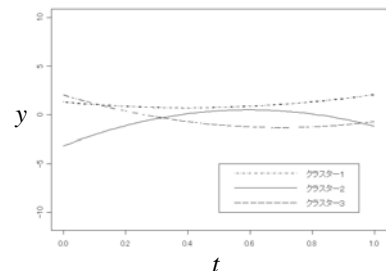
$y(t)$ の $k=3$ 個の関数主要点



各クラスターに含まれる関数データ



各クラスターの中心関数



$y(t) = a_1 + a_2 t + a_3 t^2 = b_1 + \sqrt{12}(t-1/2)b_2 + \sqrt{12}(t^2-t+1/6)b_3$ ($0 < t < 1$) の関数クラスタリングにおける局所解の個数 ($(b_1, b_2, b_3)' \sim N_3((0,0,0)', \text{diag}(\zeta_1^2, \zeta_2^2, \zeta_3^2))$) ($n=250, k=3, N=1000$)

