

Partial Distance Fuzzy C-Means

Koki Kitamori

2022 年 11 月 20 日

1 欠測値を含むデータへの FCM 法の適用

部分的距離戦略を用いて、欠測値を含むデータに対して FCM 法を適用した。
最小化する目的関数は次の式で表現される。

$$J_{fcm} = \sum_{c=1}^C \sum_{i=1}^n \sum_{j=1}^m h_{ij} u_{ci}^{\theta} (x_{ij} - b_{cj})^2 \quad (1)$$

ここで、 h_{ij} は観測値 x_{ij} の有無を表す 2 値変量で観測された場合に 1、そうでない場合に 0 の値をとる。クラスタ中心とファジィメンバシップ値の更新式は次のようになる。○はアダマール積を表す。

$$b_c = \frac{\sum_{i=1}^n u_{ci}^{\theta} (\mathbf{x}_i \circ \mathbf{h}_i)}{\sum_{i=1}^n u_{ci}^{\theta}} \quad (2)$$

$$u_{ci} = \left[\sum_{l=1}^c \left(\frac{\|\mathbf{x}_i \circ \mathbf{h}_i - \mathbf{b}_l\|^2}{\|\mathbf{x}_i \circ \mathbf{h}_i - \mathbf{b}_c\|^2} \right)^{\frac{1}{\theta-1}} \right]^{-1} \quad (3)$$

[部分的距離戦略を用いた FCM 法のアルゴリズム]

1. 要素が h_{ij} の行列 H を作る。
2. ファジィメンバシップ値 u_{ci} をランダムに定める。
3. クラスタ中心の更新とファジィメンバシップ値の更新を収束するまで繰り返す。

クラスタリング結果は以下の 2 つの指標を用いて比較した。

1. Bezdek の partition coefficient (V_{PC})

$$V_{PC} = \frac{1}{n} \sum_{i=1}^c \sum_{j=1}^n u_{ij}^2 \quad (4)$$

2. Dave の modified partition coefficient (V_{MPC})

$$V_{MPC} = 1 - \frac{c}{c-1} (1 - V_{PC}) \quad (5)$$

2 数値実験

UCI Machine Learning Repository から 3 クラス 13 変量, 178 個体からなる wine データセットを用いた。欠測値を含むデータはランダムに 20 個の観測値を欠落させて作成した。このレポートでは、欠測値を含まないデータでのクラスタリング結果と部分的距離戦略を用いた FCM 法でのクラスタリング結果の間で妥当性尺度の値にどのような差が出るかを実験した。ファジィ度は 2.0, クラスタ数は 2~6 とした。

3 実験結果

[欠測値がない場合]

	V_PC	V_MPC
C=2	0.8760833757248748	0.7521667514497496
C=3	0.7909398594655069	0.6864097891982603
C=4	0.7829819980280049	0.71064266403734
C=5	0.7440197072466368	0.680024634058296
C=6	0.7464268182671943	0.6957121819206331

[欠測値を処理した場合]

	V_PC	V_MPC
C=2	0.8760395425126468	0.7520790850252936
C=3	0.7907883399189741	0.6861825098784612
C=4	0.782683214526157	0.710244286034876
C=5	0.7434561878650163	0.6793202348312704
C=6	0.7455635445927241	0.6946762535112689

上記の実験結果から結束値をないデータでクラスタリングした場合と、欠測値 20 個含むデータに対して部分的距離戦略を用いた FCM 法でクラスタリングした場合ほとんど指標の低下がみられなかった。よってこの方法は欠測値を含む場合でのファジィクラスタリングに有効であったと考えられる。

4 ソースコード

ソースコード 1: PartialDistanceFCM.jl

```
1 using CSV, DataFrames
2 using Clustering, StatsBase, Distances
3
4 # 欠測値がない場合
5 # wine.csvは欠測値がないデータ
6 data = Matrix(CSV.read("wine.csv", header=false, DataFrame))[:, 1:end-1]
7 # クラスタ数Cを2から6としてFCM法を実行
8 # 評価指標としてBezdekのpartition coefficient (V_PC) とDaveのmodification of the V_PC(
   V_MPC)を用いた.
9 for C in 2:6
10     result = fuzzy_cmeans(data', C, 2.0, tol=1e-5, maxiter=1000)
11     w = result.weights
12     vpc = sum(w .^ 2) / length(data[:, end]) #  $V\_PC = \frac{1}{n} \sum_{i=1}^n \sum_{c=1}^C w_{ic}^2$ 
13     vmpc = 1 - (C / (C - 1)) * (1 - vpc) #  $1 - \frac{C}{C-1}(1-V\_PC)$ 
14     C == 2 && println("[欠測値がない場合]\n      |      V_PC      ,      V_MPC")
15     println("C=$C| $vpc, $vmpc")
16 end
17
18 # 欠測値がある場合
19 function update_centers!(centers, data, weights, fuzziness, H)
20     nrows, ncols = size(weights)
21     T = eltype(centers)
22     for j in 1:ncols
23         num = zeros(T, size(data, 1))
24         den = zero(T)
25         for i in 1:nrows
26              $\delta_m = weights[i, j]^{\text{fuzziness}}$ 
27             num +=  $\delta_m * (data[:, i] .* H[:, i])$ 
28             den +=  $\delta_m$ 
29         end
30         centers[:, j] = num / den
31     end
end
```

```

32 end
33
34 function update_weights!(weights, data, centers, fuzziness, dist_metric, H)
35     pow = 2.0 / (fuzziness - 1)
36     nrows, ncols = size(weights)
37     dists = pairwise(dist_metric, data .* H, centers, dims=2)
38     for i in 1:nrows
39         for j in 1:ncols
40             den = 0.0
41             for k in 1:ncols
42                 den += (dists[i, j] / dists[i, k])^pow
43             end
44             weights[i, j] = 1.0 / den
45         end
46     end
47 end
48
49 # wine2.csvは欠測値を20個含むデータ
50 data = Matrix(CSV.read("wine2.csv", header=false, DataFrame))[:, 1:end-1]
51
52 H = similar(data)
53 for j in eachindex(data)
54     ismissing(data[j]) ? H[j] = 0.0 : H[j] = 1.0
55 end
56
57 # この値を使うことはないけど適当な数値を入れとかなないとエラーになるのでなんか入れておく
58 replace!(data, missing => 99999.9)
59
60 fuzziness = 2.0
61 maxiter = 1000
62 tol = 1e-5
63
64 for C in 2:6
65     δ = Inf
66     nrows, ncols = size(data')
67     iter = 0
68
69     # 初期化
70     weights = rand(Float64, ncols, C)
71     weights ./= sum(weights, dims=2)
72     centers = zeros(nrows, C)
73     prev_centers = identity.(centers)
74
75     while iter < maxiter && δ > tol
76         update_centers!(centers, data', weights, fuzziness, H')
77         update_weights!(weights, data', centers, fuzziness, Euclidean(), H')
78         δ = maximum(colwise(Euclidean(), prev_centers, centers))
79         copyto!(prev_centers, centers)
80         iter += 1
81     end
82     #  $V_{PC} = \frac{1}{n} \sum_{i=1}^c \sum_{j=1}^n u_{ij}^2$ 
83     vpc = sum(weights.^2) / length(data[:, end])
84
85     #  $1 - \frac{c}{c-1}(1 - V_{PC})$ 
86     vmpc = 1 - (C / (C - 1)) * (1 - vpc)
87     C == 2 && println("\n[欠測値を処理した場合]\n      |      V_PC      |      V_MPC")
88     println("C=$C | $vpc, $vmpc")
89 end

```