

Title	FCM 融合型多目的クラスタリング(モデリングと最適化の理論)
Author(s)	春名, 亮; 石井, 博昭
Citation	数理解析研究所講究録 (2006), 1526: 56-60
Issue Date	2006-12
URL	http://hdl.handle.net/2433/58877
Right	
Type	Departmental Bulletin Paper
Textversion	publisher

FCM融合型多目的クラスタリング

大阪大学大学院情報科学研究科情報数理学専攻 春名 亮* Ryo HARUNA
大阪大学大学院情報科学研究科情報数理学専攻 石井 博昭 Hiroaki ISHII
Department of Information and Physical Science,
Graduate School of Information Science and Technology, Osaka University

1 はじめに

一般的に数論的意思決定問題の大半は単一の評価基準を利用することが非常に多く、それはクラスタリングについても同様の状況であるといえる。特に、クラスタリングにおいては、単一の評価基準としてデータと中心との距離の最小化を採用することは非常に典型的な方法である。従来のクラスタリングアルゴリズムは、基本的なクラスターの様々な形に従わないかもしれない唯一の評価基準が利用されている [1]。しかし、実際には本質的に複数の評価基準を持つことが多いと考えられ、一元的な価値基準で評価することができる場合のほうがむしろ稀である。そこで、クラスタリングにおいても同時に多数のクラスタリング評価関数を用いることが可能な新しいクラスタリングの方法を考えることが必要とされる。そのアプローチの一要素として、クラスターを再抽出する技法を用いて多数のクラスターの効用を評価するクラスターの適合度関数を用いている。

複数の評価基準をもつ場合として、クラスター中心と各点との偏差を利用するクラスターのコンパクト性に関する評価基準、および近傍データ同士がどの程度同じクラスターグループに属しているかを測定するデータ点の連結性に関する評価基準が用いられた。また、多目的クラスタリングでもクラスター内のデータの所属度について検討する必要があるといえる。Wang and Wu[2]は、ファジィc-平均法 (FCM : Fuzzy c-means) を2目的最適化クラスタリングに変更する方法を提案したり、M.Sato and Y.Sato[3]らは状況に依存したデータのファジィクラスタリングを多目的最適化として定式化を行い、そのパレート最適解を考察する方法を提案している。

我々は、新しい方法としてクラスターの選択方法を導入するが、0-1 整数計画問題として定式化されるので整数緩和も兼ねてクラスター選択をファジィ化する。さらにクラスター内のデータの所属にあいまいさを導入した多目的性を伴うファジィクラスタリングモデルを定式化し、パレート最適解法に依らない非線形最小化を行う方法を提案する。

2 多目的クラスタリングとファジィc-平均法

多目的クラスタリングの目標は、幾つかの異なる目的関数に対応するクラスタリングアルゴリズムを適用することによってデータ集合内でのクラスターを見つけることである。我々は1つの分割へと異なるクラスタリングアルゴリズムの出力を統合するファジィc-平均クラスタリングアルゴリズムを提案する。正確に言えば、異なるクラスタリング目的関数が与えられ、我々はデータ空間の異なる部分に対して適切な目的関数を利用して分割を

求める。

適合度関数に対する1つの可能な候補は、データ集合の再抽出のもとでクラスターの安定性があるものであり、多目的クラスタリングではしばしば発見されたクラスターが互換性のない感覚の中で矛盾する基準に遭遇する。

2.1 問題の状態

我々はデータ集合 $\mathcal{D} = \{x_1, \dots, x_m\}$ と L 個のクラスタリングアルゴリズムが与えられたと仮定して、各アルゴリズム A_i は対応する目的関数 f_i を最大にする \mathcal{D} の分割 $P^{(i)}$ へ戻すものである。正式には、

$$P^{(i)} = \{S_1^{(i)}, S_2^{(i)}, \dots, S_{C_i}^{(i)}\} = \arg \max_{P(\mathcal{D})} f_i(P(\mathcal{D})) \quad (1)$$

である。ここで、 $P(\mathcal{D})$ は \mathcal{D} の任意分割を表し、 $S_c^{(i)}$ は $P^{(i)}$ における c 番目のクラスターである。 $S \equiv \cup_i P^{(i)}$ を候補のアルゴリズムによって、発生させた全てのクラスターの集まり

$$S = \{S_1^{(1)}, \dots, S_{C_1}^{(1)}, \dots, S_1^{(L)}, \dots, S_{C_L}^{(L)}\} \quad (2)$$

である。目標は、分割 $\{P^{(i)}\}$ に基づいて「合意形成な」分割 T を見つけることである。換言すると、クラスターの集まりの目標は、 T は S におけるクラスターから得られる。

実際に、我々は全ての特徴空間に適用したクラスタリングアルゴリズムによって得られたクラスターの比較および選択をしなければならない。このため、我々はクラスタリングアルゴリズムによって使用される評価基準に外部の付加的な目的関数を導入する必要がある。目的関数を調和する方法においてクラスターの特性を測定する適合度関数 $g_j(C_i, \mathcal{D})$ とする。適合度関数 $g_j(C_i, \mathcal{D})$ を比較することによって、我々は異なるデータの部分集合(クラスター)に対する最も近似的なクラスタリング評価基準 f_i を間接的に採用することができる。

2.2 適合度関数

適合度関数 $g_j(C_i, \mathcal{D})$ はクラスター C_i と単に C_i の代わりに、全体のデータ集合 \mathcal{D} の両方に依存する。なぜなら一般的にクラスターの適合度は全てのデータ点との関連でその重要性を必要とするからである。正当な適合度関数は次の性質を持つべきである。

性質1: クラスターの評価基準 f_i または f_j を最適にするクラスタリングアルゴリズム A_j と関係を持つべきである。適合度関数 $g_j(C_i, \mathcal{D})$ の値が大きければ、さらに f_j または同等に A_j と関係がある。

性質2: 異なるクラスタリング目的関数に対して比較するべきである。換言すると、 $g_j(C_i, \mathcal{D}) > g_l(C_i, \mathcal{D})$ ならば、クラスター C_i の特性は f_l よりも f_j と関係がある。例えば、 $g_j(C_i, \mathcal{D}) = g_l(C_i, \mathcal{D})$ はクラスター C_i と C_l は同等に良く、評価基準 f_j と関係があることを意味する。

マーチンらによって、適合度関数はクラスターの安定性を基にするべきだと提案されている [1]。クラスターの安定性はデータの摂動の下でクラスタリング解法における変化を反映し、異なるクラスタリングアルゴリズムが使える。摂動は、置換の有無にかかわらず、再抽出しているデータによって生じる。安定したクラスターはたいてい選ぶに値するものである。なぜなら同じクラスターがデータ集合における軽微な変化に関係なく形成されるからである。安定したクラスターは良い分離の結果またはクラスターの緊密さをもつ。 $g_j(C_i, D)$ を計算するための擬似コードはアルゴリズム1によって与えられる。

[アルゴリズム1]

I) 以下の操作を M 回行う。

- 摂動を与えるデータ集合 D' を得るための交換の有無に関わらず D を再抽出
- 入力として A' を用いて A_j を実行して $P(D')$ を求める。
- $P(D')$ はクラスターの意味に従って D/D' におけるデータを分類することによって $P(D)$ に変換
- $\text{sim}(C_i, D)$ を計算

II) スコア l の平均を $g_j(C_i, D)$ とする。

ここで、 $\text{sim}(C_i, D)$ は任意のデータ分割 $P(D)$ を伴うクラスター C_i を比較する類似度である。

2.3 クラスターを選択

統合されたクラスター $S = \{C_1, \dots, C_M\}$ のリスト与えられ、我々は適合度関数を用いて、目標のクラスター集合 $T = \{C_1^*, \dots, C_M^*\}$ を見つける。 u はクラスターに関するメンバシップで構成されるベクトルである。ここでもしメンバシップが1ならば選択において1つのクラスターが T の中で選ばれ、そうでなければ0である。集合 T は最適なメンバシップを求めることによって構築される。

w_{ij} は C_i および C_j が集合 T 内で選ばれたならば、非対立な性質に反していることによるペナルティーを示しており、行列 W を定義する。従って u に対する非対立な性質に反している全体のペナルティーは2次項 $\frac{1}{2}u^T W u$ として考えることができる。 n_{ij} を C_i および C_j の両方に入るデータ点のメンバシップ

$$n_{ij} = \left(\sum_k q_{ki} \right) \wedge \left(\sum_k q_{kj} \right) \quad (3)$$

とする。ここで、 q_{kc} はクラスター c のデータ集合 D における x_k のメンバシップである。 w_{ij} の1つの妥当な定義は、

$$w_{ij} = \frac{n_{ij}}{\max(|C_i|, |C_j|)} = \frac{n_{ij}}{\mu} \quad (4)$$

である。直観的に言えば、 w_{ij} はまたより大きなクラスターへ割り当てられ、より小さなクラスターからのデータ点の比率を表す。 $\xi(u)$ は u において現在のクラスターによって

割り当てられずに残っているデータ点と全体の比とする。もし、完全な分布範囲のペナルティが満たされれば、 $\xi(\mathbf{u}) = 0$ である。さらに、ファジィc-平均クラスタリング法 [?] の目的関数

$$J_{\text{FCM}}^0 = \sum_k \sum_c (q_{kc})^0 \|\mathbf{x}_k - \mathbf{v}_c\|^2 \quad (5)$$

も付加する。 $\xi(\mathbf{u})$ および $\mathbf{u}^\top W \mathbf{u}$ を最小化したいのであるが、これはクラスターの適合度関数の和を最大化することと同様であるので、ファジィc-平均クラスタリング法の目的関数 [4] も追加し、3個の正の変数 $\omega_1, \omega_2, \omega_3$ を導入して、最小化目的関数を次のように定義する。

$$J = -\mathbf{s}^\top \mathbf{u} + \omega_1 \mathbf{u}^\top W \mathbf{u} + \omega_2 \xi(\mathbf{u}) + \omega_3 J_{\text{FCM}}^2 \quad (6)$$

ここで、 \mathbf{s} は適合度関数値で構成されるベクトルであり、 $s_i = g_j(C_i, \mathcal{D})$ 、 j は C_i を創生するアルゴリズム \mathcal{A}_j を記す。関数 $\xi(\mathbf{u})$ は

$$\xi(\mathbf{u}) \approx \frac{1}{m(\mathbf{d}^\top \mathbf{u} - \frac{1}{2} \mathbf{u}^\top N \mathbf{u})} \quad (7)$$

として近似され、 \mathbf{d} は C_i の大きさ d_i を要素とし、 N は n_{ij} を要素とする行列である。制約条件

$$\sum_i u_i = 1, \quad u_i \in [0, 1] \quad (8)$$

$$\sum_c q_{kc} = 1, \quad q_{kc} \in [0, 1] \quad (9)$$

の下で、(6) 式をラグランジュ未定乗数法により以下の目的関数

$$L = -\mathbf{s}^\top \mathbf{u} + \omega_1 \mathbf{u}^\top W \mathbf{u} + \omega_2 \left\{ -m \left(\mathbf{d}^\top \mathbf{u} - \frac{1}{2} \mathbf{u}^\top N \mathbf{u} \right) \right\} + \tau (\mathbf{1}^\top \mathbf{u} - 1) \\ + \omega_3 \sum_k \sum_c (q_{kc})^2 \|\mathbf{x}_k - \mathbf{v}_c\|^2 + \sum_k \eta_k \left(\sum_c q_{kc} - 1 \right) \quad (10)$$

を最小化する。ここで、 $\xi(\mathbf{u})$ を2種類に場合分けして用いることによって、メンバシップの更新式も2種類考えなければならない。ファジィc-平均法を融合した不動点反復アルゴリズム [6, 7] を示す。

Step 1: 以下に示すパラメータの初期値を与える

- C : クラスタ数
- $\omega_1, \omega_2, \omega_3$
- (7) 式における m
- 2つの小さな正数 $\varepsilon_1, \varepsilon_2$

2つのメンバシップ (u_i, q_{kc}) 、クラスター中心 \mathbf{v}_c の初期値は乱数を用いる

Step 2: クラスタ中心 \mathbf{v}_c を更新する

Step 3: 2つのメンバシップ (u_i, q_{kc}) を更新する

Step 4: もし以下の条件

$$\max |u^{NEW} - u^{OLD}| < \varepsilon_1 \quad (11)$$

$$\max |q_{kc}^{NEW} - q_{kc}^{OLD}| < \varepsilon_2 \quad (12)$$

を満足すれば終了, そうでなければ Step 2 へ戻る

3 むすび

我々は, 多目的クラスタリングおよびファジィc-平均法を融合したモデルの定式化を示した. 従来のクラスタリングの評価基準にクラスターの選択に関する評価基準を付加して, さらにクラスターの選択指標をファジィ化するとともに 0-1 整数条件を緩和したファジィc-平均法を融合した不動点反復アルゴリズムを提案した.

今後の課題として, 関数 $\xi(u)$ の近似方法の改善や新たな $\text{sim}(C_i, \mathcal{D})$ の定義などが残される.

参考文献

- [1] Martin H.C.Law, Alexander P.Topchy and Anil K.Jain : Multiobjective Data Clustering, To appear in *IEEE Computer Society on Conference on Computer Vision and Pattern Recognition*, 2004.
- [2] H.-F. Wang and G.-Y. Wu, Bi-Criteria fuzzy clustering systems, VI IFSA World Congress, Sao Paulo, Brazil, 1995, Vol.1, pp.633-636, 1995
- [3] M.Sato and Y.Sato, On a multicriteria fuzzy clustering method, Fifth IFSA World Congress, Seoul, Korea, 1993, pp.473-776, 1993
- [4] J.C.Bczdek : Pattern Recognition with Fuzzy Objective Function Algorithms, Plenum Press (1981)
- [5] Frank Hoppner, Frank Klawonn, Rudolf Kruse, Tomas Runkler, -FUZZY CLUSTER ANALYSIS AND IMAGE RECOGNITION-, WILEY (2000)
- [6] A .K. Jain and R. C. Dubes, "Algorithms for clustering data", Prentice Hall. Englewood Cliffs, New Jersey. (1988)
- [7] F. Klawonn, E.-P.Klement, Mathematical Analysis of Fuzzy Classifiers. In : X. Liu, P. Cohen, M. Berthold (eds.). *Advances in Intelligent Data Analysis*. Springer, Berlin (1997), pp.359-370