

k-means 法による Generalized Aggregation Operator を用いたクラスタリング

201011327 辻 陽介

経営工学専攻 指導教員：イリチュ（佐藤）美佳 教授

1. 研究目的

データ解析手法の一つであるクラスタリングは、個体間の類似構造に基づき、同じ特徴を持ったデータ同士に分類する手法であるが、二点間の類似性の定義は様々である。本研究で用いたクラスタリング手法の一つである k-means 法では非類似度に基づき分類を行うが、このときの非類似度には一般にユークリッド距離を用いる。

ユークリッド距離は内積を用いて表現可能であるが、この内積構造に対して Generalized Aggregation Operator(GAO)[3] を用いて、多様にある GAO の中で代表的な 3 つに対して結果の傾向を調査することを研究の目的とする。

2. Aggregation Operator

t ノルム [1],[2],[4],[5] は、三角不等式を確率分布で表現するために導入されたものである。つまり、

$$Pr\{d(p,r) \leq x+y\} \leq Pr\{d(p,q) \leq x, d(q,r) \leq y\} \quad (1)$$

(1) 式において、 $\pi_x = Pr\{d(p,q) \leq x\}$, $\pi_y = Pr\{d(q,r) \leq y\}$ とし、 $T(\pi_a, \pi_b)$ で表すと、

$$Pr\{d(p,r) \leq a+b\} \leq T(\pi_a, \pi_b) \quad (2)$$

という関係式が得られるが、(2) 式における T が満たすべき性質を取り出したのが、t ノルムである。

t ノルムは以下のように定義される。

$[0, 1] \times [0, 1] \rightarrow [0, 1]$ からなる演算 が以下の条件を満たすとき t-ノルムと呼ばれる。

$$(T.1) \quad T(a, b) = T(b, a) \quad (\text{対称性})$$

$$(T.2) \quad T(a, c) \leq T(b, d)$$

whenever $a \leq b, c \leq d$ (単調性)

$$(T.3) \quad T(a, 1) = T(1, a) = a$$

$T(a, 0) = T(0, a) = 0$ (境界条件)

$$(T.4) \quad T(T(a, b), c) = T(a, T(b, c)) \quad (\text{結合性})$$

Aggregation Operator(AO) は t ノルムの 4 つの条件のうち、3 つの条件 (対称性、単調性、境界条件) を満足するものである。

3. Generalized Aggregation Operator(GAO)

Generalized Aggregation Operator(GAO)[3] は、AO を線形空間上の直積空間上で定義した演算である。

GAO は以下のように定義される。

線形空間を X とすると、 $X \times X \rightarrow [0, 1]$ からなる演算 $\tilde{\rho}$ が以下の条件を満たすとき、Generalized Aggregation Operator と呼ばれる。

$$(\tilde{A}.1) \quad \tilde{\rho}(\mathbf{a}, \mathbf{b}) = \tilde{\rho}(\mathbf{b}, \mathbf{a}) \quad (\text{対称性})$$

$$(\tilde{A}.2) \quad \tilde{\rho}(\mathbf{a}, \mathbf{c}) \leq \tilde{\rho}(\mathbf{b}, \mathbf{d})$$

whenever $\mathbf{a} \leq \mathbf{b}, \mathbf{c} \leq \mathbf{d}$ (単調性)

$$(\tilde{A}.3) \quad \tilde{\rho}(\mathbf{a}, \mathbf{1}) = \tilde{\rho}(\mathbf{1}, \mathbf{a}) = \alpha \quad \alpha \in [0, 1]$$
$$\tilde{\rho}(\mathbf{a}, \mathbf{0}) = \tilde{\rho}(\mathbf{0}, \mathbf{a}) = 0 \quad (\text{境界条件})$$

本研究で用いる GAO は以下の三つである。

Generalized Algebraic Product(GAP):

$$\tilde{\rho}(\mathbf{a}, \mathbf{b}) = \mathbf{ab}^t$$

Generalized Hamacher Product(GHP):

$$\tilde{\rho}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{ab}^t}{\mathbf{a1}^t + \mathbf{b1} - \mathbf{ab}^t}$$

Generalized Einstein Product(GEP):

$$\tilde{\rho}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{ab}^t}{2K - (\mathbf{a1}^t + \mathbf{b1} - \mathbf{ab}^t)}$$

また、GAO はいずれも以下の条件を満たしている。

$$\sum_{k=1}^K a_k = \sum_{k=1}^K b_k = \sum_{k=1}^K c_k = \sum_{k=1}^K d_k = 1, \quad 2 \leq K \leq n \quad (3)$$

4. GAO を用いたクラスタリング

k-means 法では非類似度にユークリッド距離を用いるが、ユークリッド距離は、内積を用いて表現できる。本研究では内積に GAO を用いたクラスタリングを行う。ユークリッド距離を GAO で表現すると、以下のようになる。

$$d(\mathbf{x}_i, \bar{\mathbf{x}}_c) = \{\tilde{\rho}(\mathbf{x}_i, \mathbf{x}_i) - 2\tilde{\rho}(\mathbf{x}_i, \bar{\mathbf{x}}_c) + \tilde{\rho}(\bar{\mathbf{x}}_c, \bar{\mathbf{x}}_c)\}^{\frac{1}{2}} \quad (4)$$

ここで、 \mathbf{x}_i は個体 i を表すベクトル、 $\bar{\mathbf{x}}_c$ はクラスター c の重心を表している。(4) 式に対して、の GAO(GAP,GHP,GEP) を用いて、分類の結果を調査した。

5. 数値例

数値実験では人工データを使用した。まず、1番目から25番目までの2変数のデータを平均(4,6)、分散1の正規乱数で発生させ、これをクラスター1とした。次に、クラスター2として26番目から50番目までの2変数データを平均(8,2)、分散0.8の正規乱数により作成した。この時1番目から50番目までのデータを $\mathbf{x}_1, \dots, \mathbf{x}_{50}$ とし、 i 番目($i=1, \dots, 50$)のデータを $\mathbf{x}_i = (x_{i1}, x_{i2})$ と表すとき、各データが変数について和が1であるという条件を満たすように次の変換を行う。

$$\tilde{\mathbf{x}}_i = (\tilde{x}_{i1}, \tilde{x}_{i2}), \quad \tilde{x}_{ia} \equiv \frac{x_{ia}}{\sum_{b=1}^2 x_{ib}}, \quad a = 1, 2 \quad (5)$$

以上の変換を行ったデータに対して、GAOを用いたクラスタリングで解析すると、それぞれ図1、図2、図3の結果となった。この時一つのGAOに対する図の数は重心が変化しなくなるまでの更新回数を示している。

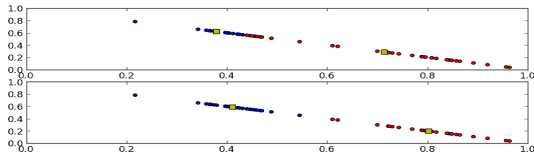


図1 GAPによるクラスタリング

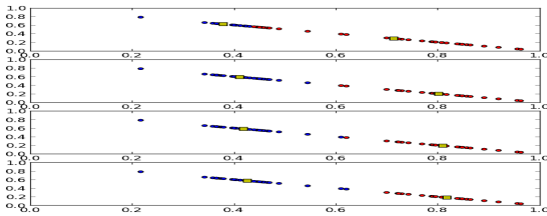


図2 GHPによるクラスタリング

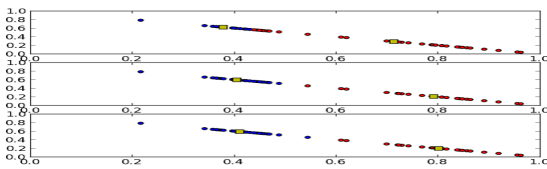


図3 GEPによるクラスタリング

図1、図2、図3より、GHP、GEP、GAPの順番に重心の更新回数が多いことが分かり、さらにGHPはほかの二つとは違う分類の結果を示した。一方GEPはGAPとは最終的に同じ分類の結果を示した。

次に、1番目から1000番目までの2変数のデータを平均(2.5,7.5)、分散0.8の正規乱数で発生させ、これをクラスター1とした。また、クラスター2として1001番目から2000番目までの2変数データを平均(7.4,2.6)、分散0.7の正規乱数により作成した。初め的人工データと同様、(5)式による処理を行い、GAOを用いたクラスタリングを行った。結果は図4のように得られ、すべて同じ結果が得られた。

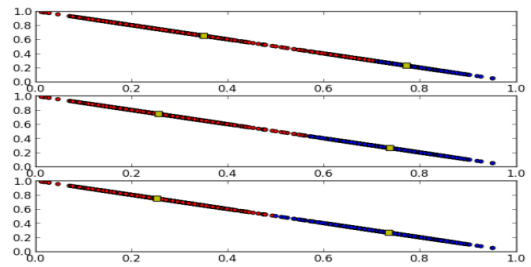


図4 GAOによるクラスタリング

6. 結論

本論文ではGAOを用いたクラスタリングについての研究を行った。k-means法で用いられている非類似度で最も一般的なユークリッド距離を内積で表現し、GAOを用いることでそれぞれのGAO(GAP、GEP、GHP)の分類の相違性を調査したものである。

本研究の数値実験より、分類構造が明確なデータの場合は、GAP、GEP、GHPによる結果に違いが生じるが、分類構造が不明確なデータの場合は、違いが生じないということがわかった。

参考文献

- [1] E. P. Klement., R. Mesiar., E. Pap., *Triangular Norms*, Kluwer Academic, 2000
- [2] K. Menger., *Statistical Metrics*, Proc. Nat. Acad. Sci., USA, vol.28, pp. 535-537, 1942.
- [3] M. Sato-Ilic, Generalized Aggregation Operator Based Nonlinear Fuzzy Clustering Model, *Engineering Systems through Artificial Neural Networks*, vol. 20, pp. 493-500, 2010.
- [4] B. Schweizer., A. Sklar., *Probablistic Metric Space*, Dover Publications, 2005
- [5] 中島信之, t-ノルム, t-コノルム通覧(1)-概説と歴史, 日本ファジィ学会誌 Vol. 11, No. 4, pp. 561-576, 1999