

近隣集合に基づく特徴変数の類似性評価に関する研究

前野 良隆 市野 学

東京電機大学 大学院 理工学研究科 情報社会学専攻

1. はじめに

シンボリック・データ・アナリシスは、データから有益な情報を獲得する統計手法の一般化を目指している。既存の多くの統計手法は、対象の記述が量的特徴に基づくことを前提としている。したがって、量的特徴でも区間を値とする場合や質的な特徴が混在する場合には、既存の統計手法の直接的な一般化は難しい。このような場合の手法として数量化理論が知られているが、量的特徴を区間に分けるという任意性の残る操作が必要であったり、数多くのダミー変数を用いることによる特徴の増大など、適用上対処すべき問題点を残している。

本研究では、カルテシアン・システム・モデルと呼ばれる数学モデルを用いる。この数学モデルでは、各対象が量的特徴と質的特徴の混在する形式で記述されることを認め、しかも特徴の型を統一することを目的とした特別な処理は必要としない。このカルテシアン・システム・モデルを基に、対象が記述される多次元空間において、対象間の相対的關係を記述する近隣集合の概念を定義する。次に近隣集合を用いて、多次元空間に対象群が構成する「幾何学的に薄い構造」を検出するための方法について考察する。また、提案手法の有用性を実験によって確認する。

2. 相対近隣と近隣集合

2-1. 相対近隣 (Nearest Relative Neighborhood)

量的・質的データを扱う事が出来る計算モデルとしてカルテシアン・システム・モデル(以下CSM) [2]があり、相対近隣はCSMに基づく。CSMでは、領域を定義するカルテシアン・ジョインという演算があり、サンプル対のジョインを作成した時、その領域内に含まれる他のサンプルの数 p を Generality p であるという。

d 個の特徴 $F_x(x=1,2,\dots,d)$, n 個のサンプルの有限集合 $= \{ i_1, i_2, \dots, i_n \}$ によって記述されているとする。

特徴集合 F の p 番目の特徴 F_p において、任意のサンプル対 i_p, j_p によって形成される閉区間を想定する。その閉区間に包含されるサンプルを k_p とし、 k_p の集合 i_{jp} を、

$$i_{jp} = \{ k_p | \min(i_p, j_p) \leq k_p \leq \max(i_p, j_p) \} \quad (2-1)$$

とする。ただし、 $\min(i_p, j_p), \max(i_p, j_p)$ は、それぞれ i_p と j_p の小さい方の値、 i_p と j_p の大きい方の値をとる演算である。

$$\text{generality}(i_p, j_p | F_p) = | i_{jp} | \quad (2-2)$$

ここで、 i_{jp} に含まれるサンプルの数を以下のように定義する。ただし、 $| \cdot |$ は集合 \cdot の基数を表す。

注目する特徴 F_p について、互いに隣接しているサンプル対、つまり、(2-2)式 $\text{generality}(i_p, j_p | F_p)$ の値が k となるサンプル対を、特徴 F_p について k 相対近隣 (k -Nearest Relative Neighborhood) であると定義する。特に $\text{generality } k$ が 0 の時を単に相対近隣と言う。

2-2. 入れ子構造

入れ子構造とは、あるサンプル対の中に別のサンプル対が入った構造のことをいう。単調構造は最大サンプルと最小サンプルの中に内在する Generality が 1 個ずつ減少していく入れ子構造を有している。単調構造の量的特徴の例を図 2.1, 質的特徴の例を表 2.2 に示す。

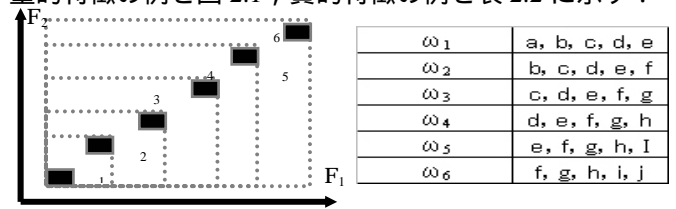


図 2.1 単調構造(量的特徴)

表 2.2 単調構造(質的特徴)

2-3. 近隣集合 (Neighborhood set)

近隣集合は相対近隣に基づく。

近隣集合とは、ある基準となるサンプルの相対近隣の関係にあるサンプル i_{n-1}, i_{n+1} と、基準となったサンプル i_n の集合であると定義する。

$$n(i_n | F) = \{ i_{n-1}, i_n, i_{n+1} \} \quad (2-3)$$

例として図 2.1 を用いて説明する。最初に 1 次元について考えたい。まず F_1 についてサンプル集合の近隣を見ていくと、下記ようになる。

$$\begin{aligned} n(\omega_1 | F_1) &= \{ \omega_1, \omega_2 \} & n(\omega_2 | F_1) &= \{ \omega_1, \omega_2, \omega_3 \} \\ n(\omega_3 | F_1) &= \{ \omega_2, \omega_3, \omega_4 \} & n(\omega_4 | F_1) &= \{ \omega_3, \omega_4, \omega_5 \} \\ n(\omega_5 | F_1) &= \{ \omega_4, \omega_5, \omega_6 \} & n(\omega_6 | F_1) &= \{ \omega_5, \omega_6 \} \end{aligned}$$

F_1 について見ると最初と最後が 2 つ、その間が自分を含む前後の 3 つずつになっているのがわかる。 F_2 , および、 F_1, F_2 の場合も同様の結果が得られる。このように、単調な構造とは全ての特徴において各サンプルの近隣集合が同じになるという特性を有する。

3. k 相対近隣を考慮に入れた単調性の程度の評価尺度

単調性の程度を評価する尺度として本提案手法は以下の式(3-1)で与える。この手法は近隣集合および近隣集合を $\text{generality } k$ まで増やすことを許可した k 相対近隣に基づく。

A study on the similarity of feature variables based on neighborhood sets
Yoshitaka MAENO, Manabu ICHINO
Department of Information and Arts, Tokyo Denki University

$$M = 1 - \left(\frac{1}{2 \sum_{i=0}^k (2+k+i) + (3+2k)(n-2(k+1)) \sum_{k=1}^n |NRN_p(\omega_k|F)|} \right) \times \frac{2k}{n-2}$$

$$0 \leq M \leq 1 \quad \dots (3-1)$$

ここで k は k 相対近隣を示し、評価値は 1 をもっとも良い評価値とし、0~1 の間に推移する。

4. 性能評価実験

4-1. ノイズへの耐性

提案手法を評価するため、線形関数、指数関数、対数関数に逐次的にノイズを加えたデータを用いた。

各種関数を $v=f(u)$ によって記述されているとし、 X 、 Y をそれぞれ独立な標準正規乱数とする。

線形構造では説明変数 u に標準正規乱数 X を発生させた時、ピアソンの積率相関係数を用いて目的変数 v を以下の式(4-1)で与える。

$$v = pX + \sqrt{1-p^2}Y \quad (4-1)$$

v と u のピアソンの積率相関係数の値は p となる。

指数関数、対数関数では、説明変数 u に一様乱数 Z [-5,5] を発生させた時、目的変数 v を以下の式(4-2)で与える。

$$v = f(u) + \alpha Y \quad (4-2)$$

ここでは、ノイズの大きさを指定するパラメータであり、は 0.0 から 0.1 刻みで 1.0 まで変化させたものである。

実験に使用したデータはいずれも、サンプル数は 100、特徴数は 12 である。特徴 1 を u とおき、特徴 2~12 はを 0.0 から 1.0 まで変化させたときの v である。

ノイズがない状態でのピアソンの積率相関係数との比較を図 4.1、ノイズを加えていったデータに対しての実験結果を図 4.2 に示す。

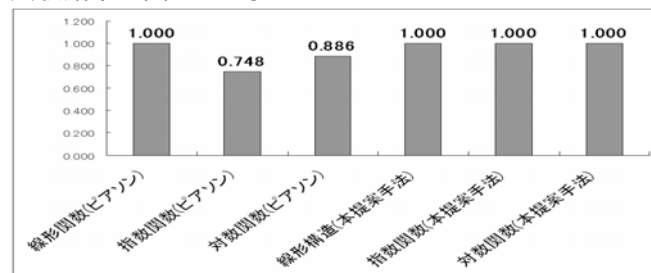


図 4.1 ピアソンの積率相関係数との比較

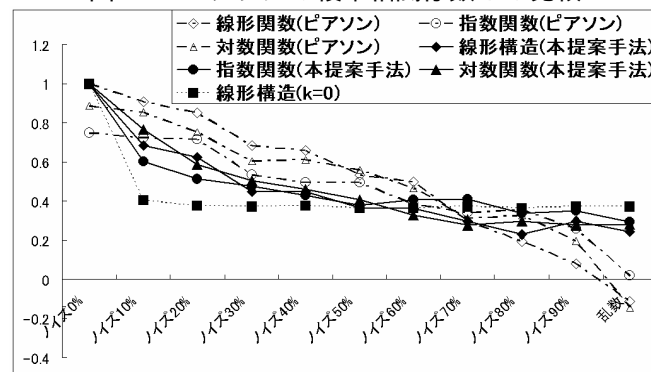


図 4.2 ノイズへの耐性

k 相対近隣を考慮しない ($k=0$) 時、ノイズのないデ

ータには単調構造を抜き出すことが可能である。しかし、ノイズが少しでもある場合、単調構造を抜き出すことが不可能である。

次に k 相対近隣を考慮に入れた本提案手法ではノイズが加わったデータに対してもノイズの量によって評価値が下がっており、ピアソンの積率相関係数に近い評価値を示すことができた。これはノイズの量に対して、本提案手法が適性に評価ができたと考える。

また指数関数、対数関数に関して、ピアソンの積率相関係数を用いた場合、ノイズのない場合でも約 0.750、0.900 の値であったのに対し、本提案手法ではノイズのない場合、 k 相対近隣を考慮することによって 1.000 の値をとることが示された。またノイズに対しても、ノイズの量によって評価値が下がる傾向にあり、ノイズの量に依存した合理的な評価値が得られている。

4-2. 質的な特徴を含む人工データ

提案手法を評価するため、質的な特徴を含む特徴集合によって形成される人工データを用いて評価実験を行った。実験に使用したデータは、サンプル数は 12、特徴数は 5 であり、質的な特徴とする。特徴 4、5 を質的な特徴の単調構造、残りの特徴は規則性を持たない冗長な特徴である。

本提案手法は CSM を用いているため、質的なデータであっても、単調な構造は抜き出すことが可能である。本提案手法の実験結果を図 4.3 に示す。

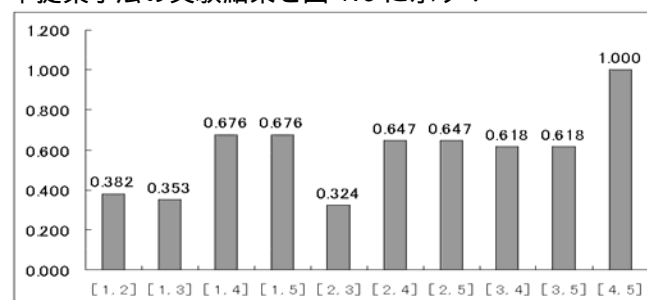


図 4.3 実験結果

実験の結果、本提案手法は質的な特徴が内在しているデータであっても評価することができた。また冗長な特徴が含まれたデータに対しても、単調な構造を抜き出すことが可能である。この結果、量的特徴、質的な特徴の混在したシンボリックデータに対応していることが実証できた。

5. まとめ

本研究では特徴間の類似性の尺度を見るために近隣集合を提案した。また k 相対近隣を用いることにより、 k 相対近隣を考慮に入れた単調性の程度の評価尺度を定義した。簡単な人工データによる性能評価実験において、提案手法の有用性を示した。

参考文献

- [1] M.Ichino, Detection of Monotonic Chain Structures in Mixed Feature Type Multidimensional Data, 2005.
- [2] 市野学, 矢口博之, 野中武志, “幾何学的厚みに基づく相関係数”, 電子情報通信学会論文誌 A, Vol.J85-A, No.4, pp.490-494, (2002).