

論理最小化アルゴリズムによる確率分布からの命題抽出の検討

3W-1

吉澤 有美 三浦 勤 稲積宏誠

青山学院大学理工学部経営工学科

1.はじめに

本論文では、離散確率分布からの命題抽出という教師なし帰納学習を検討する。教師なし機能学習とは、分類クラスなしに事例が与えられたときにその事例から適切な命題などを獲得することである。すでに情報理論に基づく、確率分布からの命題抽出[1]が提案されているが、命題の簡単化と確率分布との関係、さらに多値属性を持つ命題表現を扱ったものは少ない。そこで本論文では、命題表現の簡単化に際して、多値属性にも適用可能でヒューリスティックアプローチを含む論理最小化アルゴリズムである MINI[2]を用いて、より有効な命題表現を生成する方法を検討する。

その手順は以下のとおりである。

1. 与えられた確率分布から KL(Kullback-Leibler) 情報量を最小にする古典確率ベクトルと論理ベクトルを求める。
2. 確率分布から論理ベクトルに対応した重みベクトルを求める。
3. MINI のアルゴリズムに重みベクトルを用いて命題表現の簡単化を行う。

2. 命題表現

属性値数 p_i の属性 X_i の属性値を $x_i^j (j = 1, 2, \dots, p_i)$ とし、属性 X_i をベクトル $\{x_i^1, x_i^2, \dots, x_i^{p_i}\}$ と表す。ただし、

$$x_i^j = \begin{cases} 1 & ; \text{if 属性値として } x_i^j \text{ を選ぶ;} \\ 0 & ; \text{otherwise;} \end{cases}$$

である。その結果、各属性は $\{0,1\}$ ベクトルとして表現され、これをパートと呼ぶ。また、 n 属性からなる項は、 $\wedge_{i=1}^n X_i = \{(x_1^1, x_1^2, \dots, x_1^{p_1}), \dots, (x_n^1, x_n^2, \dots, x_n^{p_n})\}$ と表し、長さ $\sum_{i=1}^n p_i$ の $\{0,1\}$ ベクトルとして表現され、キューブと呼ぶ、ここで、最小項は次のように表される。

$$\phi_j = \wedge_{i=1}^n X_i^{k(i)}; \quad k(i) \in \{1, 2, \dots, p_i\}, \quad j = \{1, 2, \dots, \prod_{i=1}^n p_i\}$$

ただし、 $X_i^{k(i)} = \{x_i^1, \dots, x_i^{k(i)}, \dots, x_i^{p_i}\}$;

$$x_i^a = \begin{cases} 1 & ; \alpha = k(i); \quad k(i) \in \{1, 2, \dots, p_i\}; \\ 0 & ; \alpha \neq k(i); \quad k(i) \notin \{1, 2, \dots, p_i\}; \end{cases}$$

その結果、任意の命題は最小項 ϕ_j を用いて、

$$\vee_{j=1}^{\prod_{i=1}^n p_i} a_j \phi_j; \quad a_j \in \{0,1\}$$

と表現できる。ただし、 j は ϕ_j を 2 進数とみなしたときの降順に順序づけるものとする。このようにして求められた係数 a_j のベクトル表現を論理ベクトルと呼ぶ。

例えば、3 つの属性値をとる属性 X_1 が存在したとする。この属性 X_1 が x_1^2 という属性値をとる場合、このパートは 010 である。また、2 変数（属性 X_1 (3 値 x_1^1, x_1^2, x_1^3) と属性 X_2 (2 値 x_2^1, x_2^2)）からなる命題 $(x_1^1 \vee x_1^3) \wedge x_2^2$ は、通常キューブで 10101 と表現する。この命題は最小項表現で、 $x_1^1 x_2^1 \vee x_1^3 x_2^2 = (100 \ 01) \vee (001 \ 01)$ と表現できる。このとき係数 a_j は

$$a_2 = a_6 = 1, \quad a_1 = a_3 = a_4 = a_5 = 0,$$

となり、この命題の論理ベクトルは、 $<0,1,0,0,0,1>$ と表すことができる。

3. 確率分布からの論理ベクトル抽出

無差別原理を用いて確率分布から古典確率ベクトルと論理ベクトルを求める方法を示す[1]。

無差別原理は複数の事象のどれが生起するかに関して情報がない時、最小項表現されている各事象に対して、生起確率は等しいと考える原理である。ここで 1 変数（2 値）のトートロジーを考える。これは、 $x_1^1 \vee x_1^2$ と書くことができ、ベクトル表現は $<1,1>$ となる。一方、トートロジーが持つ情報量は 0 であり、情報量 0 の確率分布は $(1/2, 1/2)$ となる。従って、論理ベクトル $<1,1>$ と古典確率ベクトル $(1/2, 1/2)$ は対応する。同様にして、多値属性をもつ 2 変数（2 属性と 3 属性）の場合も、情報量 0 に相当する論理ベクトル $<1,1,1,1,1,>$ は古典確率ベクトル $(1/6, 1/6, 1/6, 1/6, 1/6, 1/6)$ に対応することになる。以上により、多値属性を想定した場合にも論理ベクトル成分の 1 の数を m とした時には、古典確率ベクトルの成分を $1/m$ か 0 で表すことになる。

古典確率ベクトルを真の確率分布 $p = \{p_1, \dots, p_m = 1/m, p_{m+1}, \dots, p_N = 0\}$ と仮定し、与えられた確率分布をモデル $q = \{q_1, \dots, q_N\}$ と見なし、式(1)に示す KL 情報量を最小にする m を求めることによって論理ベクトルを決定する。

$$I(p; q) = \sum_{i=1}^N p_i \log p_i - \sum_{i=1}^N p_i \log q_i \\ = \log \frac{1}{m} - \frac{1}{m} (\log q_1 + \dots + \log q_m) \quad \cdots (1)$$

さらに、ここで得られた論理ベクトルの要素 1 に対する確率分布の値から構成されるベクトルを重みベクトルとして定義する。

5. 命題の簡単化

多値属性を持つ論理を簡単化できるアルゴリズムに MINI[2]がある。このアルゴリズムはキューブの順序、パートの順序にヒューリスティックを含んでおり、その順序により論理表現（キューブ数、キューブの形）が異なる。MINI ではキューブ数を最小にすることを目的にこのヒューリスティックを用いているが、本研究では確率分布を命題の簡単化に影響させるヒューリスティックを提案する。以下にアルゴリズムを示す。

1. 論理ベクトル成分が 1 の事例をキューブ表現 ($\{0, 1\}$ ベクトル) で表しこのリストを $F = v(cube f_i)$ とする。各々のキューブ f_i のキューブ重みを、その事例の重みベクトル要素の値とする。

2. disjoint sharp 演算(以下④)*

$U \oplus F$ によってキューブから成る \bar{F} を生成、さらに $U \oplus \bar{F}$ によって異なる構造の F を生成し、キューブの数をある程度減少させる。この演算は先に演算をするキューブ、またその中の後に演算するパートほど、マージされやすいという傾向がある。

従って、キューブ重みの降順に演算を行う。演算後のキューブ重みは、そのキューブを構成している最小項のキューブ重みの和として表現される。

以下キューブの構造に応じてキューブ重みは更新されるものとする。

3. expand 演算**

④演算で得られた \bar{F} に対して F のキューブを拡張する。拡張されたキューブ(以下 $expand(f_i)$)は F の範囲内で拡張するため、 F に含まれる他のキューブを除く。

* < disjoint sharp process >

$A = \pi_1, \pi_2, \dots, \pi_p, B = \mu_1, \mu_2, \dots, \mu_p$ とすると

$$A \oplus B = C = v'_i, C_i$$

$$C_i = (\pi_1 \wedge \mu_1), (\pi_2 \wedge \mu_2), \dots, (\pi_{i-1} \wedge \mu_{i-1}), (\pi_i \wedge \mu_i), \pi_{i+1}, \dots, \pi_p$$

$$U = Universe, F = v(cube f_j)$$

$$U \oplus (U \oplus F)$$

で F を生成します。

一ブあるいはその一部をカバーする。そこで、 $expand(f_i)$ にカバーされているキューブを F から取り除き、 F のキューブ数をカウントする。この拡張演算におけるキューブの順序に更新されたキューブ重みを用いる。

4. reshape 演算***

パートの 2箇所が異なるペアを F から見つけ出し、最小項が影響を受けないような他の disjoint なキューブのペアに変える。

5. キューブ数に変化がなくなるまで 3,4 を繰り返す。

本アルゴリズムをアンケート結果からの命題抽出の例題[1]に適用し、その性質を検討した。提案アルゴリズムでは、常にキューブ重みとして確率分布の情報をヒューリスティックとして活用しているため、得られた論理表現に対する説明機能を付加できると考えられる。

6. おわりに

確率分布からの命題抽出に関する研究には、命題抽出の過程に確率分布を用いているが、命題を簡単化する際には、確率分布を考慮していないものが多い。本研究では確率分布を命題の簡単化に影響させるヒューリスティックを示すとともに、確率分布からの命題抽出を多値属性を持つ命題表現に拡張した。

参考文献：

[1] 森田千絵 月本洋：最尤法と無差別原理を用いた確率分布からの帰納学習、人工知能学会誌、Vol. 10, No. 4, pp297-304(1997)

[2] S. J. Hong, R. G. Cain, D. L. Ostapko: MINI:A Heuristic Approach for Logic Minimization, IBM J. Research and Development, pp. 443-458(1974).

* < expand process >

$cube f = \pi_1, \pi_2, \dots, \pi_p$ を $\bar{F} = v(cube g_i)$ に反して expand する

$$H(f; k) = \{g_i | f \text{ and } g_i \text{ are } k - \text{conjugate}\}$$

ビットごとの OR をとり、 k 番目のパートだけ

$NULL(0)$ となるとき f と g_i は $k - \text{conjugate}$ である。

$$Z(f; k) = \text{すべての } H(f; k) \text{ のビットごとの OR}$$

$$SPE(f; k) = \pi_1, \pi_2, \dots, \pi_{k-1}, \overline{Z(f; k)}, \pi_{k+1}, \dots, \pi_p$$

$$expand(f) = SPE(\dots SPE(SPE(f; \sigma(1)); \sigma(2)); \dots; \sigma(p))$$

*** < reshaping process >

F に含まれる 2つのキューブ $f_a = \pi_1, \pi_2, \dots, \pi_p, f_b = \mu_1, \mu_2, \dots, \mu_p$, f_a と f_b が $part i$ と j で異なり、

π_i が μ_i をカバーし、 π_j が μ_j をカバーしないとき

$$f'_a = \pi_1, \pi_2, \dots, (\pi_j \wedge \overline{\mu_j}), \dots, \pi_p, f'_b = \mu_1, \mu_2, \dots, (\pi_i \vee \mu_i), \dots, \mu_p,$$