

入学年度 平成 10 年度

学籍番号 10117950

氏名 長谷川 友治

論文題目 属性間のクラス条件付き相関に強い決定木生成手法

犬塚研究室

## 1 はじめに

決定木学習は、決定木という表現形式が人間にも理解しやすく、大規模データにも適用できるアルゴリズムが多数存在することから、データマイニングの手法としても利用されている。しかし、クラスに対する属性間の相関によっては、C4.5のような一つの属性のみで分割した時に得られる重要度を用いた属性選択では、本当に重要な属性が選択できない問題がある。本論文では属性間のクラス条件付き相関に強い属性選択基準を提案する。

## 2 属性間のクラス条件付き相関

一般的に、全事例集合  $D$  上で成り立つ属性間の相関は冗長な属性が含まれていることになる。それとは別に、あるクラス  $c$  に所属する事例の部分集合  $D_c$  のみに現れる相関を属性間のクラス条件付き相関と呼び、その共起パターンを以下のように表現する。

$$(a_i = v_i) \xrightarrow{D_c} (a_j = v_j)$$

$v_i, v_j$  はそれぞれ  $a_i, a_j$  の取りうる属性値のひとつである。このクラス条件付き相関は、クラス  $c$  に所属する事例の中だけで発生している相関であり、この相関がクラス  $c$  についての性質を説明していると考えられる。

しかし、複数の属性の相関によってクラスを決定付ける場合、ID3 や C4.5 では単一の属性の評価しかしないため見逃してしまうことがある。

## 3 相関に強い属性選択基準の提案

いくつかの属性の間には何らかの相関がある問題を考える。ここで、ある属性  $a_i$  以外のすべての属性の値を一定にしたまま、属性  $a_i$  の値の変化によるクラスの変化があれば、 $a_i$  がクラスに影響を与えていることが分かる。このとき、 $a_i$  と  $a_i$  以外の一定にした値の間には何らかの相関がある。このような  $a_i$  を選ぶことで、相関を考慮した属性選択ができると考えられる。そこで  $a_i$  に関する相関がどの程度あるかを、 $a_i$  以外の属性の値の組合せ  $V_{a_i}$  が生じたときのクラスの条件付きエントロピーとして量的に計算し、重要度として用いる。

$$\operatorname{argmax}_{a_i \in A} - \sum_{V_{a_i}} \frac{|X_{V_{a_i}}|}{|X|} \sum_{j=1}^l \frac{|X_{V_{a_i}}^{c_j}|}{|X_{V_{a_i}}|} \log_2 \frac{|X_{V_{a_i}}^{c_j}|}{|X_{V_{a_i}}|}$$

ここで、 $X_{V_{a_i}}$  は  $V_{a_i}$  と同じ値 ( $a_i$  の値を除く) を取る事例の集合を表し、 $X_{V_{a_i}}^c (\subseteq X_{V_{a_i}})$  は  $X_{V_{a_i}}$  内の事例でクラスが  $c$  である事例の集合を意味している。

## 4 実装と実験結果

上の基準式を用いた実装にはハッシュ表を用いることで、漸近的な計算量が C4.5 と同じになる実装を行なった。

実験として、相関性のある 11-multiplexor 問題と UCI のデータベースを用いて C4.5, MINITREE, 提案手法で比較実験を行なった。評価方法は、テストデータが提供されている monk1, monk3 はそれを用いて評価し、それ以外の問題は 10-fold cross validation を行ない評価した。

結果を表 1 に、得られた木の例を図 1, 2 に示す。

表 1: 実験結果 Size:木のサイズ,E:誤分類率

DB	C4.5		MINITREE		提案手法	
	Size	E	Size	E	Size	E
mux11-0	204.8	0.2	<b>35.8</b>	<b>0.0</b>	46.6	<b>0.0</b>
mux11-5	246.2	15.8	78.4	14.9	<b>47.0</b>	<b>14.7</b>
monk1	<b>18.0</b>	24.3	38.0	<b>5.6</b>	28.0	8.3
monk3	<b>12.0</b>	2.8	<b>12.0</b>	2.8	17.0	<b>0.0</b>
mux6	22.6	36.2	<b>17.4</b>	2.9	21.8	<b>0.0</b>
lenses	<b>6.4</b>	<b>18.3</b>	6.9	21.7	<b>6.4</b>	<b>18.3</b>
car	<b>168.6</b>	7.4	189.1	13.5	182.0	<b>6.8</b>
nursery	503.6	<b>2.8</b>	634.8	5.2	<b>495.6</b>	<b>2.8</b>
zoo	18.2	<b>0.6</b>	<b>1.0</b>	59.8	27.3	16.9
vote	<b>13.5</b>	<b>4.4</b>	18.7	5.3	19.6	14.0
tic-tac-toe	125.5	<b>14.6</b>	113.5	24.4	<b>28.6</b>	29.6
mushroom	<b>33.1</b>	<b>0.0</b>	675.5	0.1	134.7	7.3

mux11- $x$  は 5000 事例の 11-multiplexor 問題に  $x\%$  noise を加えたもの

C4.5 と比較して、相関のある mux11, mux6, monk1 などでは決定木の性能の向上が見られた。mux の問題では、提案手法によって生成された決定木の方が非常にコンパクトになって知識としてより分かりやすい表現になっている。monk1 の問題では、分類精度が非常に向上した。それ以外の問題でも、事例数が十分ある car, nursery, lenses などは C4.5 とほぼ同等な結果が得られた。

ノイズがある mux11-5 では、MINITREE と比較して木が小さく精度も勝っており、ノイズに強いと言える。

しかし、事例空間に対して事例数が少ない問題や、zoo のような各事例に ID を付けている属性および冗長性のある属性を含む問題では、 $X_{V_{a_i}}$  の要素数が少なく重要度が十分に計算できない。そのため、zoo, vote, tic-tac-toe, mushroom では C4.5 と比べて決定木の性能が悪くなった。

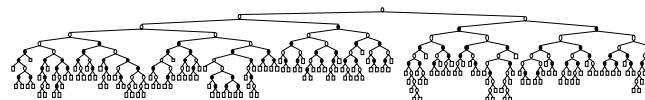


図 1: C4.5 によって生成された木 (mux11-5)



図 2: 提案手法によって生成された木 (mux11-5)

## 5 まとめ

本論文では、決定木生成における相関に強い属性選択基準の提案を行なった。この基準を用いた実験では、属性間に相関のある問題について性能の高い決定木が生成でき、ノイズが含まれている場合にも耐性があった。また、属性間に明確な相関がない問題でも訓練事例の数が十分ある問題においては、C4.5 と同等か良い結果が得られた。ただし、事例数が少ない場合や、冗長な属性、事例に ID を与えるような属性が含まれる問題については、重要度が十分に計算できない。

今後の課題は、学習を阻害する属性の削除、および欠損値、連続値の対応が挙げられる。

## 参考文献

- [1] 長谷川, 松井, 犬塚, 世木: 決定木生成における属性間の相関に強い属性選択基準, 2002 年度人工知能学会全国大会 (発表予定)