

帰納学習によるファジィ決定木の生成

正員 櫻井 茂 明 (東 芝)

非会員 荒 木 大 (東 芝)

Generating a Fuzzy Decision Tree by Inductive Learning

Shigeaki Sakurai, Member, Dai Araki, Non-member (Toshiba Corporation)

ID3 algorithm can automatically acquire a decision tree from a set of training examples. However, ID3 can only deal with the distinct data values. This paper presents a fuzzy decision tree which expresses some fuzzy classification rules, and an algorithm to induce a fuzzy decision tree from the training examples including numerical or fuzzy data values. This new algorithm, called IDF, has a labeling procedure in each decision tree expanding step to make effective fuzzy classification items. These items are used to do fuzzy decisions in the branch nodes. Fuzzy decision tree can also give some classification results with certainty ratios. The authors examined and exemplified the efficiency of this algorithm by some numerical experiments.

キーワード：ID3 アルゴリズム, ファジィ決定木, IDF アルゴリズム

1. はじめに

近年の人工知能の研究において、分類規則を自動的に獲得する帰納学習の技術が広く研究され、その応用が進みつつある。その一つとして、複数の属性と、それに対して与えられた分類クラスを訓練事例として、多数の訓練事例から属性と分類クラスの間的一般的な規則を見つけ出す手法が提案されている。

このような手法には大きく分けて二つの手法があり、一つは、ID3⁽¹⁾に代表される決定木形式の分類規則を生成するアルゴリズムであり、もう一つは、AQ11⁽²⁾に代表される選言的ブール形式の分類規則を生成するアルゴリズムである。

本論文では、これらの手法のうち決定木形式の分類規則を生成する手法に着目する。

オリジナルのID3アルゴリズムには、訓練事例を表現する属性は、離散的な値しかとれないという適用上の制限がある。従って、数値や人間が主観的に与えたデータを含んだ属性をID3で取り扱うためには、それらのデータを離散的なデータに変換する必要がある。数値で与えられるような連続的な属性の扱いについては、分類能力の高い境界値を発見する区間分割ア

ルゴリズムを組み込んだINDECT⁽³⁾アルゴリズムを提案した。しかしながら、通常数値で与えられるようなデータには、雑音が含まれているので、発見した境界値は一応の目安にすぎず、境界値付近の値に対して判断誤りを起こしやすいという欠点があった。また、「寒い」、「暑い」などの人間が主観的に与えたデータを、ID3アルゴリズムでは、互いに独立した離散的なデータとみなして取り扱っている。しかしながら、本来このような主観的なデータは境界を明確に決めることはできないので、独立したものとみなすことはできない。

本論文では、あいまいなデータを取り扱う有効な手法であるファジィ集合理論⁽⁴⁾⁽⁵⁾を適用することにより、決定木の表現形式を拡張する。更に、あいまい性を加味した決定木（以下、ファジィ決定木と呼ぶ）を生成するアルゴリズムを提案する。このアルゴリズムは、属性の取り得る値をあいまいに分類するファジィ集合を動的に生成しながら、事例集合のファジィ分割を行うことにより、ファジィ決定木形式の分類規則を生成する。生成されたファジィ決定木による推論では、分岐ノードであいまいな判断を行い、末端ノードに到達した評価対象の確信度の和で分類クラスの判定

を行う。従って、従来の決定木のような択一的な判定は行われず、確信度のついた分類結果が得られる。

本論文を以下のように構成する。第2章では、ファジィ決定木とその構成方法について説明する。第3章では、提案した方法の有効性を示すべく行った数値実験の方法と実験結果について説明する。

2. ファジィ決定木とファジィ決定木生成アルゴリズム

本章では、ファジィ集合理論⁽⁴⁾⁽⁵⁾により、表現形式を拡張したファジィ決定木について説明する。また、あいまい性を含んだ訓練事例からファジィ決定木を学習する IDF アルゴリズムについて説明する。

〈2・1〉 ファジィ決定木 ファジィ決定木は、属性値のペアで表現されている評価対象から分類クラスを確信度つきで判定する規則を表現している。このファジィ決定木では、評価対象の属性値にファジィな値を用いることができる。

ファジィ決定木は、属性の取り得る値をあいまいに分類するファジィ集合〔以下、ファジィ分岐判断項目

(FCI) と呼ぶ〕を分岐ノードにもち、確信度のついた分類クラスを末端ノードにもつ。ここで、確信度は分類クラスに該当する確信の度合を 0.0~1.0 の数値で表した値である。

ファジィ決定木を用いた推論では、分岐ノードのもつ FCI により、評価対象があいまいに判断され、複数の下位のノードに伝搬する。末端ノードは、到達した評価対象の確信度と分類クラスの確信度を掛けて部分解を生成する。最終的には、分類クラスごとに部分解を合計することにより分類結果を得る。

例えば、図1のファジィ決定木で、「HUMIDITY が 80% であり、TEMPERATURE が more or less hot である」という評価対象の分類クラスを判定することを考える。評価対象に与えられる初期確信度を 1.0 とし、最初に属性値「more or less hot」を属性「TEMPERATURE」で判断する。すなわち、「TEMPERATURE」のもつ FCI により、確信度 0.2 $[=1.0(0.25/0.25+1.0)]$ で「mild」と、確信度 0.8 $[=1.0(1.0/0.25+1.0)]$ で「hot」とあいまいに判断する。評価対象の伝搬した分岐ノードに対して同様な

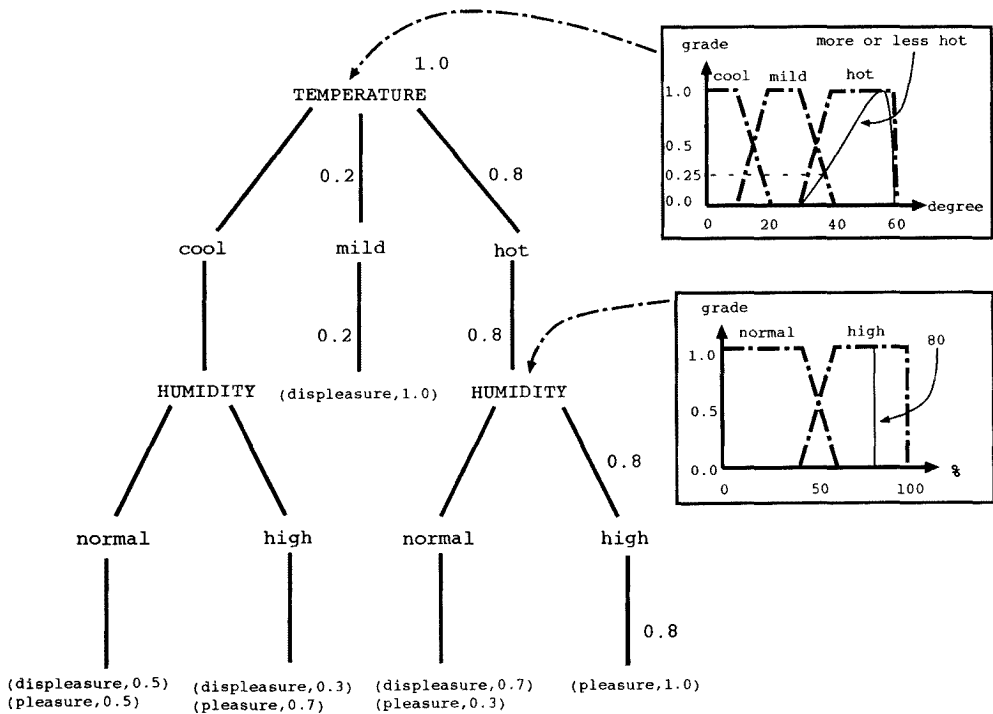


図1 ファジィ決定木
Fig. 1. Fuzzy decision tree.

判断を繰り返していくことにより、図1に示すように確信度が伝搬する。末端ノードごとに生成される部分解から「(pleasure, 0.8), (displeasure, 0.2)」という分類結果を得る。

IDF アルゴリズムは、このような推論を行うファジィ決定木を、属性値と分類クラスの既知な訓練事例集合から生成するアルゴリズムである。

〈2・2〉 IDF アルゴリズム IDF アルゴリズム⁽¹⁾は n_A 個の属性と n_C 種類の分類クラスからなる多数の訓練事例を入力として、決定木形式の分類規則を生成する帰納学習アルゴリズムである。ID3は相互情報量が最大となる属性を選択し、選択した属性の属性値に従って事例集合を分割し、決定木を成長させる。この決定木の成長を、分割された事例集合に含まれる訓練事例の分類クラスが一つに決まるまで行う。

IDF アルゴリズムは ID3 アルゴリズムを基本とし、与えられる訓練事例の属性値としてファジィ集合を取り扱うために、FCIの生成ステップ[ステップ3(a)]を新たに追加し、相互情報量の計算方法[ステップ3(b)]および事例集合の分割方法[ステップ3(c)]を改良している。以下の項において、これらの部分について詳しく説明する。

(1) ファジィ分岐判断項目の生成 この項では、ファジィ決定木の構成要素となるファジィ分岐判断項目 (FCI) の生成方法について説明する。

一般に、分類クラスが同一の訓練事例群ごとに属性値の代表値の平均値を求めたならば、その平均値の周りに分類クラスが等しい属性値が集まっている。また、分散の小さいものほど平均値の周りに属性値が分布し、大きなものほど平均値から離れて属性値が分布

している。従って、平均値を中心として、分散が小さいならばメンバシップ関数の広がり小さくし、分散が大きければメンバシップ関数の広がりを大きくするように、FCIのメンバシップ関数を生成すれば分類クラスの判別を行うという観点で、有効な FCI が生成できると考えられる。

この考えを基にして、訓練事例から FCI を生成するアルゴリズムを表2のように構成する。このアルゴリズムでは、計算を簡単にするため、属性 A_i に対して生成する n_i 個の FCI がそれぞれもつメンバシップ関数 $m_{fik}(x)$, ($1 \leq k \leq n_i$) を(1)式の形式とする。以下では、この台形型のメンバシップ関数を $(\alpha_{ik}, \beta_{ik}, \gamma_{ik}, \delta_{ik})$ と表す。

$$m_{fik}(x) = \begin{cases} 0, & x \leq \alpha_{ik} \\ \frac{x - \alpha_{ik}}{\beta_{ik} - \alpha_{ik}}, & \alpha_{ik} < x < \beta_{ik} \\ 1, & \beta_{ik} \leq x \leq \gamma_{ik} \\ \frac{x - \delta_{ik}}{\gamma_{ik} - \delta_{ik}}, & \gamma_{ik} < x < \delta_{ik} \\ 0, & \delta_{ik} \leq x \end{cases} \quad \dots\dots\dots (1)$$

表2のアルゴリズムで導入しているしきい値 T_1 は、確信度の和が小さくなった事例集合の分割を中止することにより FCI の数を抑制している。一方、しきい値 T_2 は、分離度が高くなった事例集合の分割を中止することにより FCI の数を抑制している。

次に、表2のアルゴリズムで行われている、各種の計算方法について説明する。以下の説明において、 $C = \{c_1, c_2, \dots, c_n\}$ を分類クラスの集合、 S を分割の対象にしている事例集合、 S_{c_k} を S に含まれる訓練事例

表 1 IDF アルゴリズム
Table 1. IDF algorithm.

1. すべての訓練事例を一まとめにして対応付けたノードを生成する。
2. ノードに属する訓練事例のうちで、同一の分類クラスをもつ訓練事例の確信度の和の割合がしきい値以上ならば、そのノードを末端ノードとし、確信度のついた分類クラスをラベル付けする。しきい値より小さければ、そのノードを分岐ノードとする。
3. 分岐ノードに対して、
 - (a) 各属性に対して、FCI を生成する。
 - (b) 各属性の相互情報量を計算し、最大の値を取る属性 A_k を、その分岐ノードにおける判断属性とする。
 - (c) 判断属性 A_k の FCI_{fka} , ($k=1, 2, \dots, n_k$) によって、分岐ノードに属する訓練事例を n_k 個のファジィ部分集合に分割する。
 - (d) 分割された個々のファジィ部分集合に対して、新しいノードを生成する。このとき元のノードと新たに生成したノードを結ぶリンクには、対応する FCI をラベル付ける。
 - (e) 新しく生成したノードに対して、[ステップ2]からの手続きを再帰的に適用する。

表 2 ファジィ分岐判断項目生成
アルゴリズム
Table 2. Algorithm for making FCI.

1. ノードに割り当てられている事例集合から分割を開始する。
2. 事例集合に含まれる各々の訓練事例に対して、属性 A_i の属性値の代表値を計算する。
3. 分類クラスが同一の訓練事例群ごとに、確信度を重みとした属性値の代表値の平均値と分散を計算する。
4. 分類クラスを平均値の小さい順に並べる。
5. 分類クラスごとに、属性 A_i の FCI を生成する。
6. FCI に基づいて事例集合をファジィ分割し、確信度の和が設定するしきい値 T_1 より大きなファジィ部分集合を更に分割を行う候補とする。
7. 分割を行う候補から、分類クラスに対する分離度が最小となるファジィ部分集合を選択する。
8. 選択するファジィ部分集合がないか、選択したファジィ部分集合の分離度がしきい値 T_2 より大きければ分割を終了する。さもなければ、[ステップ7]で選択したファジィ部分集合を新たな事例集合として、[ステップ2]からの処理を繰り返す。

のうち分類クラスが $c_k \in C$ となる訓練事例の集合、事例集合 S, S_{c_k} に含まれる訓練事例の確信度の和を $|S|, |S_{c_k}|$ とする。また、 S の j 番目の訓練事例を t_j , S の i 番目の属性を A_i , t_j の A_i に関する属性値を v_{ij} , t_j の確信度を p_j , v_{ij} のもつメンバシップ関数を $m_{v_{ij}}(x)$ とする。ただし、 x は属性の取り得る値を表す。

ステップ 2 で計算する属性値 v_{ij} の代表値 r_{ij} は (2) 式により計算する。

$$r_{ij} = \frac{\int_{m_{v_{ij}}(x) > 0} x m_{v_{ij}}(x) dx}{\int_{m_{v_{ij}}(x) > 0} m_{v_{ij}}(x) dx} \quad (2)$$

この値はメンバシップ関数 $m_{v_{ij}}(x)$ と x 軸により囲まれた領域の重心の x 座標に相当している。

ステップ 3 で計算する分類クラス c_k に関する代表値の平均値は、(3) 式により計算する。

$$E_{c_k} = \frac{1}{|S_{c_k}|} \sum_{t_j \in S_{c_k}} r_{ij} p_j \quad (3)$$

また、分類クラス c_k に関する代表値の分散は (4) 式により計算する。

$$V_{c_k} = \frac{1}{|S_{c_k}|} \left\{ \sum_{t_j \in S_{c_k}} (E_{c_k} - r_{ij})^2 p_j \right\}^{1/2} \quad (4)$$

ステップ 5 では属性 A_i の FCI を生成する。代表値の平均値でソートされた隣接する分類クラスを c_{k1} , $c_{k2} \in C$, 分類クラス c_{k1} , c_{k2} に対応する平均値を $E_{c_{k1}}$, $E_{c_{k2}}$, 分散を $V_{c_{k1}}$, $V_{c_{k2}}$ とする。このとき、隣接する FCI f_{ik1} , f_{ik2} のメンバシップ関数の境界は (5), (6) 式により計算する。

$$\gamma_{ik1} = \alpha_{ik2} = E_{c_{k1}} + \frac{V_{c_{k1}}}{2(V_{c_{k1}} + V_{c_{k2}})} (E_{c_{k2}} - E_{c_{k1}}) \quad (5)$$

$$\delta_{ik1} = \beta_{ik2} = E_{c_{k2}} - \frac{V_{c_{k2}}}{2(V_{c_{k1}} + V_{c_{k2}})} (E_{c_{k2}} - E_{c_{k1}}) \quad (6)$$

すなわち、(5), (6) 式の値が台形型のメンバシップ関数 $m_{f_{ik1}}(x)$, $m_{f_{ik2}}(x)$ の x 座標となる。

$$m_{f_{ik1}}(x) = (\alpha_{ik1}, \beta_{ik1}, \gamma_{ik1}, \delta_{ik1})$$

$$m_{f_{ik2}}(x) = (\alpha_{ik2}, \beta_{ik2}, \gamma_{ik2}, \delta_{ik2})$$

ステップ 7 で計算する分離度は (7) 式により計算する。

$$u(S) = \max_{c_k \in C} \frac{|S_{c_k}|}{|S|} \quad (7)$$

この値は事例集合 S に含まれる訓練事例に関して、確信度の和が最も大きな分類クラスの占有率を表す。この値が 1 ならば事例集合に含まれる訓練事例はすべて同じ分類クラスをもつことになる。

表 3 事例集合

Table 3. Training examples.

	A_1	...	A_i	...	A_{n_a}	C	p
t_1	v_{11}	...	v_{i1}	...	v_{n_a1}	c_1	1.0
t_2	v_{12}	...	v_{i2}	...	v_{n_a2}	c_1	0.9
t_3	v_{13}	...	v_{i3}	...	v_{n_a3}	c_2	1.0
t_4	v_{14}	...	v_{i4}	...	v_{n_a4}	c_2	0.9
t_5	v_{15}	...	v_{i5}	...	v_{n_a5}	c_2	0.8
t_6	v_{16}	...	v_{i6}	...	v_{n_a6}	c_2	0.8
t_7	v_{17}	...	v_{i7}	...	v_{n_a7}	c_1	0.8
t_8	v_{18}	...	v_{i8}	...	v_{n_a8}	c_1	0.8

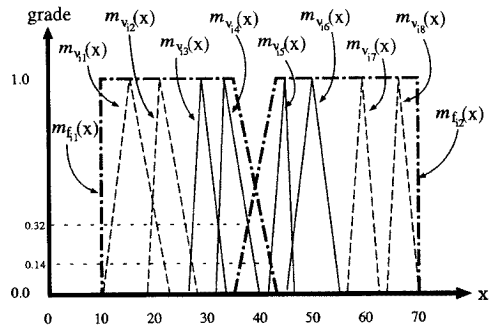


図 2 v_{ij}, f_{i1}, f_{i2} のもつメンバシップ関数
Fig. 2. Membership functions for v_{ij}, f_{i1}, f_{i2} .

例として、表 3 の事例集合に対して属性 A_i に関する FCI を生成する過程を説明する。それぞれの属性値 v_{ij} は図 2 に示すメンバシップ関数をもつ。また、実線のメンバシップ関数が分類クラス c_1 に対応し、破線のメンバシップ関数が分類クラス c_2 に対応する。

最初、アルゴリズムにより、表 3 の事例集合に対して、図 2 に一点鎖線で示すメンバシップ関数 $m_{f_{i1}}(x)$, $m_{f_{i2}}(x)$ をもつ FCI f_{i1}, f_{i2} が生成される。図 2 からわかるように、 f_{i1}, f_{i2} に割り当てられている事例集合は、分類クラスの判別という観点でまだ十分に分割されていない。従って、 f_{i1}, f_{i2} に割り当てられている事例集合を更にファジィ分割する。

その結果、図 3 に一点鎖線で示すメンバシップ関数 $m_{f_{i12}}(x)$, $m_{f_{i11}}(x)$, $m_{f_{i21}}(x)$, $m_{f_{i22}}(x)$ をもつ FCI $f_{i12}, f_{i11}, f_{i21}, f_{i22}$ が生成される。各 FCI に割り当てられている事例集合は分類クラスに関して十分に分割されているので、アルゴリズムを終了する。このとき、各 FCI には表 4 に示す事例集合が割り当てられている。ただし、同表では、属性 A_i 以外の属性の記述を省略している。

(2) 相互情報量の計算 IDF アルゴリズムは、確信度のついた訓練事例を取り扱っている。従って、

表 4 FCI f_{i12} , f_{i11} , f_{i21} , f_{i22} に割り当てられている事例集合

Table 4. Training examples with FCI f_{i12} , f_{i11} , f_{i21} , f_{i22} .

A_i	C	p	A_i	C	p	A_i	C	p	A_i	C	p
t_1	v_{11}	c_1	1.0	t_2	v_{12}	c_1	0.24 $\left(=0.9 \times \frac{0.37}{1.0+0.37}\right)$	t_4	v_{14}	c_2	0.22
t_2	v_{12}	c_1	0.66	t_3	v_{13}	c_2	0.81 $\left(=1.0 \times \frac{1.0}{2.4+1.0}\right)$	t_5	v_{15}	c_2	0.70
t_3	v_{13}	c_2	0.19	t_4	v_{14}	c_2	0.68 $\left(=0.9 \times \frac{1.0}{1.0+0.32}\right)$	t_6	v_{16}	c_2	0.8
				t_5	v_{15}	c_2	0.10 $\left(=0.8 \times \frac{0.14}{0.14+1.0}\right)$				

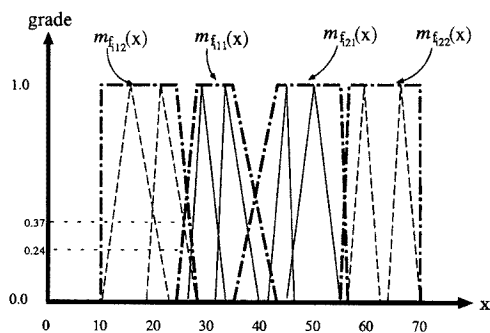


図 3 FCI のもつメンバシップ関数
Fig. 3. Membership functions for FCI.

ID3 アルゴリズムのように訓練事例の要素数に基づいて相互情報量⁽⁶⁾を計算したのでは、確信度の異なる訓練事例を一律に取り扱うことになる。一方、確信度の大きな訓練事例は、確信度の小さな訓練事例に比べて情報量が多いと考えられる。従って、確信度を考慮して相互情報量を計算する必要がある。そこで、表 1 のステップ 3(b) では、訓練事例の確信度の和に基づいて相互情報量を計算する。この計算において、属性 A_i の FCI f_{ik} に割り当てられる訓練事例 t_j の確信度は、FCI f_{ik} に対する属性値 v_{ij} の一致度 $g_{ik}(t_j)$ [(8) 式] で、確信度を正規化することにより計算される。

$$g_{ik}(t_j) = \max_x \{ \min(m_{f_{ik}}(x), m_{v_{ij}}(x)) \} \quad \dots\dots\dots (8)$$

(3) 事例集合の分割 IDF アルゴリズムでは、訓練事例の属性値がファジィ集合として与えられる。一般に、与えられるファジィ集合には種々のものがあり、類似したファジィ集合や包含関係にあるファジィ集合を同時に与える可能性もある。従って、訓練事例として与えられた属性値をそのまま用いた分割では、類似したファジィ集合や包含関係にあるファジィ集合

までも異なったものとして取り扱うことになり、決定木の枝の数が多くなり過ぎる。そこで、表 1 のステップ 3(c) では属性を有効に分割する FCI により、事例集合を分割する。この FCI を用いた分割により、事例集合を適当な数にファジィ分割することができる。

3. 数値実験と評価

本章では、提案した IDF アルゴリズムに対して行った ID3 アルゴリズムおよび INDECTS アルゴリズムとの比較実験について説明する。また、実験結果から IDF アルゴリズムの性能を評価する。

〈3・1〉 実験 1: 数値データからの学習 実験 1 では、数値データを属性値にもつ訓練事例から分類規則を学習する。実験用のサンプルとして、表 5 に示す値を平均とし、分散 1 をもった数値データを考える。ただし、ID3 では数値データをそのまま扱うことはできないので、表 5 に与えられている各属性値の中心値を境界とする区間をあらかじめ設定し、数値データを離散データに変換した事例を用いる。すなわち、属性 A_1 に対して、 $a_1 < 4.5$, $4.5 \leq a_1 < 7.5$, $7.5 \leq a_1$ という三つの区間を設定する。従って、属性 A_1 の属性値が 5.0 ならば、2 番目の区間 $4.5 \leq a_1 < 7.5$ に変換される。

この実験サンプルは、 A_1 と A_2 , A_3 と A_4 , A_5 と A_6 をそれぞれペアとして同一の分布をもつ属性が 3 通りあるので、本来ならば、二つの属性を判定するだけで正しい分類結果が得られる。しかしながら、事例のもつ雑音の影響を避けるために、三つ以上の属性値を評価して分類性能の劣化を吸収する必要がある。実験方法としては、正規乱数を用いて、各事例の出現頻度が同じになるように 100 個の訓練事例と 100 個の評価事例を発生させる。次に、100 個の訓練事例で決定木を学習し、学習された決定木を用いて 100 個の評価事例の分類クラスを判定する。この分類クラスと、本来与えられている評価事例の分類クラスとを比較し、決定木の性能評価を行った。

表 5 数値データに対する性能評価サンプル
Table 5. Evaluation examples for numerical data.

A_1	A_2	A_3	A_4	A_5	A_6	C	p
3	3	3	3	3	3	c_1	1
9	6	9	6	9	6	c_2	1
3	6	3	6	3	6	c_3	1
6	3	6	3	6	3	c_4	1
6	6	6	6	6	6	c_5	1

表 6 数値データに対する実験結果
Table 6. Experiment results for numerical data.

	正解率	第二候補までの正解率
ID 3	85.0 %	—
INDECTS	80.8 %	—
IDF	87.6 %	96.6 %

ID 3, INDECTS, IDF の各アルゴリズムで学習した決定木に対する正解率を表 6 に示す。ここで、表 6 の正解率は、5 通りの訓練事例と評価事例を組み合わせた実験で得られた正解率の平均値である。また、第二候補までの正解率とは、IDF が確信度のついた分類結果を出力することに着目して、2 番目に大きい確信度を与える分類結果までに正解となった場合の正解率である。表 6 からわかるように、IDF の正解率は INDECTS の正解率よりも高くなっている。従って、境界付近のデータに関して判断誤りを起こしやすかった INDECTS の欠点を、あいまいな判断を行うことにより回避できたとわかる。また、IDF の正解率は区間をあらかじめ設定して学習した ID 3 よりも高くなっている。従って、雑音を考慮した境界を生成することにより、ある程度の雑音を含んだ事例も正しく判定できたとわかる。更に、第二候補まで含んだ正解率はほとんど 100% となっている。従って、分類結果に付与された確信度が妥当な値となっていたとわかる。

〈3・2〉 実験 2：幅をもったデータからの学習

実験 2 では、 $[a, b]$ という幅をもった値を属性値にもつ訓練事例から分類規則を学習し、各アルゴリズムの性能を比較する。実験用のサンプルとして、表 7 に示す幅をもったデータを用いる。ただし、実際に与えられるデータは、表 7 に示す下限と上限を平均とし、分散 $\sqrt{2}$ をもった値を幅の下限と上限とするデータである。INDECTS では幅をもったデータを取り扱うことはできないので、幅をもったデータの中心値をあらかじめ計算して、数値データに変換した事例を用いる。また ID 3 では、実験 1 と同様な方法でこの中心

表 7 幅をもったデータに対する性能評価
サンプル

Table 7. Evaluation examples for interval data.

A_1	A_2	A_3	A_4	A_5	A_6	C	p
3	3	(2.5, 3.5)	(2.5 3.5)	(2, 4)	(2, 4)	c_1	1
9	6	(8.5, 9.5)	(5.5 6.5)	(8, 10)	(5, 7)	c_2	1
3	6	(2.5, 3.5)	(5.5 6.5)	(2, 4)	(5, 7)	c_3	1
6	3	(5.5, 6.5)	(2.5 3.5)	(5, 7)	(2, 4)	c_4	1
6	6	(5.5, 6.5)	(5.5 5.5)	(5, 7)	(5, 7)	c_5	1

表 8 幅をもったデータに対する実験結果
Table 8. Experiment results for interval data.

	正解率	第二候補までの正解率
ID 3	85.8 %	—
INDECTS	82.2 %	—
IDF	88.6 %	98.0 %

値を更に離散データに変換した事例を用いる。

実験方法としては、正規乱数を用いて生成した事例に対して、実験 1 と同様な方法で、決定木の性能評価を行った。ID 3, INDECTS, IDF の各アルゴリズムで学習した決定木に対する正解率を表 8 に示す。同表からわかるように、やはり IDF の正解率が一番高くなっている。従って、幅をもったデータに関しても、IDF は属性を有効に分割しながら、正解率の高い分類規則を生成するとわかる。

〈3・3〉 検 討

(1) ファジィ決定木の分類能力 ファジィ決定木で数値データやあいまい性を含んだデータを取り扱った場合に、境界値近辺の値に対して高い分類性能を示す傾向がある。この理由として、数値データやあいまい性を含んだデータが与えられた場合に、分岐ノードであいまいな判断しか行われないので、上位の分岐ノードで誤った判断がなされたとしても、下位の分岐ノードの判断でこれを回避する効果が得られたからと考えられる。

更に、ファジィ決定木は確信度のついた分類結果を得ることができるので、結果がどの程度信頼できるか数値により判定することができる。

(2) 訓練事例の記述力 IDF アルゴリズムでは、数値で与えられた属性値とファジィ集合で与えられた属性値を同時に取り扱うことができる。すなわち、ある訓練事例の属性値として、「10」、「[8, 12]」、「およそ 7」といった属性値が同時に与えられたとしても、(9)式に示すメンバシップ関数をそれぞれ定義することにより、各属性値を IDF 内で一律に取り扱

うことができる。

$$\left. \begin{array}{l} m_{10}(x) = (10, 10, 10, 10) \\ m_{(8,12)}(x) = (8, 8, 12, 12) \\ m_{\text{およそ}}(x) = (6, 7, 7, 8) \end{array} \right\} \dots\dots\dots (9)$$

また、IDF では個々の訓練事例がもつあいまい性を、初期確信度として与えて学習に反映することができる。すなわち、分類クラスが c_1 となる割合が 0.8、分類クラスが c_2 となる割合が 0.2 となる訓練事例が与えられたとするならば、確信度 0.8 で分類クラス c_1 、確信度 0.2 で分類クラス c_2 とした二つの訓練事例を用いることにより学習を行うことができる。このように、訓練事例の属性値の表現方法の幅が広がったことが IDF の大きな利点である。

4. おわりに

数値データやあいまい性を含んだデータからあいまい性を加味した分類規則を学習する IDF アルゴリズムを提案し、IDF の有効性を数値実験により確認した。

IDF を適用することにより、数値やあいまい性を含んだ訓練事例からも、分類規則を生成することができる。従って、帰納学習を適用できる問題が広がったと考えられる。特に、診断を行うような問題で、不確実なデータを用いて推論を行う必要がある場合に、IDF が適用可能と考えられる。

今後の課題としては、現行の IDF では、分類クラスが同一な訓練事例群ごとに属性値の代表値の平均値と分散を計算して、FCI がもつメンバシップ関数を生成している。しかしながら、代表値をベースにした考え方では、属性値のメンバシップ関数の形状を十分に反映しているとは限らない。従って、代表値によらないメンバシップ関数の生成方法を検討する必要がある。また、FCI の生成アルゴリズムでは、FCI を分割する操作しか行っていないので、与えられる事例集合によっては、FCI の数が多くなり過ぎてしまう。従って、統合操作などをアルゴリズムに組み込んで FCI の数が多くなり過ぎないようにする必要がある。更に、IDF を実際問題に適用してその有効性を確認していきたい。

日ごろより御指導をいただいている(株)東芝研究開発センター 西島誠一副所長、河野 毅部長、小島昌一課長に感謝いたします。

(平成 4 年 10 月 23 日受付、同 5 年 2 月 15 日再受付)

文 献

- (1) J. R. Quinlan: "Induction of Decision Trees", *Machine Learning*, 1, 71 (1985)
- (2) R. S. Michalski & R. L. Chilausky: "Learning by being told and learning from examples", *Int. J. Policy Analysis and Information Systems*, 4, No. 2, 125 (1980)
- (3) 荒木・小島: 「数値データによる決定木の帰納学習」, 人工知能学誌, 7, No. 6, 992 (平 4)
- (4) L. A. Zadeh: "Fuzzy Sets", *Information Control*, 8, 338 (1965)
- (5) L. A. Zadeh: "Fuzzy set as a basis for a theory of possibility", *Fuzzy Sets & Systems*, 1, No. 1, 3 (1978)
- (6) 西田・竹田: 「ファジィ集合とその応用」, 森北出版
- (7) 櫻井・荒木: 「ファジィ理論を適用した知識獲得」, 第 15 回計測自動制御学会知能システムシンポジウム資料, p. 169 (平 4)
- (8) 櫻井・荒木: 「あいまい性を含んだ訓練事例からの学習」, 情報処理, 92-AI-84, 31 (平 4)



櫻 井 茂 明 (正員)

平成元年東京理科大学理学部応用数学科卒業。3 年同大学大学院理学研究科数学専攻修士課程修了。同年、(株)東芝入社。現在、同社研究開発センターに勤務。機械学習の研究および開発に従事。情報処理学会会員。



荒 木 大 (非会員)

昭和 61 年大阪大学工学部通信工学科卒業。63 年同大学大学院工学研究科博士前期課程修了。同年、(株)東芝入社。現在、同社研究開発センターに勤務。知識獲得支援技術、機械学習の研究および開発に従事。情報処理学会、人工知能学会、IEEE 会員。