

子育てに関するアンケートからの 医療機関選択ルールの抽出と同定

十倉伸太郎, 村田忠彦, 松原光也



文部科学省私立大学社会連携研究推進拠点
関西大学政策グリッドコンピューティング実験センター

Policy Grid Computing Laboratory,
Kansai University
Suita, Osaka 564-8680 Japan
URL : <https://www.pglab.kansai-u.ac.jp/>
e-mail : pglab@jm.kansai-u.ac.jp
tel. 06-6368-1228
fax. 06-6330-3304

関西大学政策グリッドコンピューティング実験センターからのお願い

本ディスカッションペーパーシリーズを転載、引用、参照されたい場合には、ご面倒ですが、弊センター（pglab@jm.kansai-u.ac.jp）宛にご連絡いただきますようお願い申し上げます。

Attention from Policy Grid Computing Laboratory, Kansai University

Please reprint, cite or quote WITH consulting Kansai University Policy Grid Computing Laboratory (pglab@jm.kansai-u.ac.jp).

子育てに関するアンケートからの医療機関選択ルールの抽出と同定

十倉伸太郎¹, 村田忠彦^{2,3}, 松原光也^{1,3}

Extraction and Identification of Selection of Medical Institutions Rules from the Questionnaire about bringing up a child

Shintaro Tokura¹, Tadahiko Murata^{2,3}, Mitsuya Matsubara^{1,3}

概要

近年、マルチエージェントシミュレーション (MAS) が様々な分野から注目されている。MAS を用いた社会シミュレーションにおいて、人間の複雑な意思決定をモデル化することが重要な課題となっている。そこで本研究では、医療機関選択に関するアンケートと GIS の座標データを用いて、エージェントの行動ルールとするための人間の意思決定を表す If-Then ルールを生成する。If-Then ルールの生成において、モデル設計者によって前件部の属性を固定化したルール生成と J.R. Quinlan によって提案された ID3 による自動ルール生成の識別率を比較する。さらに、その結果を踏まえて ID3 をアンケートからの If-Then ルールの抽出に用いた際に必要な改良点を示し、また前件部に用いる属性値のグルーピングを行うことでルール数の減少とアンケートデータの分析を行った。

Abstract

Recently, multi-agent simulation (MAS) attracts attentions from various. Modeling human decision-making is a challenge in social simulation using MAS. In this research, we generate if-then rules represent human decision-making for agent's action rules using questionnaire about bringing up a child and coordinate data of GIS. In generating if-then rules, we compare the distinction ratio with fixed number of attributes in conditional part to those generated by ID3 proposed by Quinlan. Furthermore, we show a point to be improved in extracting if-then rules from questionnaire using ID3, and a way to decrease the number of rules by grouping attribute values in the conditional part of if-then rules.

キーワード : MAS, ID3, GIS

Keywords: Multi-agent simulation, ID3, Geographic information system

1 関西大学大学院

2 関西大学総合情報学部

3 関西大学政策グリッドコンピューティング実験センター

1. はじめに

政策提案の背景には統計手法による分析や数値実験によって裏付けされたデータが必要となる。近年、マルチエージェントシミュレーション (MAS) を用いた社会シミュレーションが多様な分野の研究者から注目されている。MAS により社会現象をモデル化してシミュレーションする事で、現実的に時間やコストの制限から得る事が困難な知見や考察を短期間で得ることが可能である。

MAS を用いた社会シミュレーションにおいてエージェントの意思決定プロセスを人間のそれに基づいてモデル化し、行動ルールを生成することは重要な課題である。人間の意思決定は大変複雑であり、モデル作成者の主観に基づくエージェントの意思決定のモデル化は信頼性に欠けると考えられる。この課題を解決するためには社会学や心理学からの人間の意思決定に関する知見やその知識が求められると共に、アンケートなどの現実に基づいたデータを用いて、その属性を解析することで意思決定プロセスをモデル化することが必要となる。データから知識を獲得する手法は多く提案されている。このとき、意思決定プロセスの表現方法を考える必要があるが、構造の理解や知識の解釈が容易であることから決定木による If-Then ルールの抽出が多く用いられている。決定木構築法の代表的なものに、J.R. Quinlan[1]によって提案された ID3 がある。その基本原理は獲得情報量の期待値を最大とする属性を選択することで決定木を構築するものである。これは情報量の導入により、決定木の構築において最も分類効率の良い入力属性を選択することができる[2,3,4]。

本研究では、大阪府吹田市内の 6 つの幼稚園のいずれかに属する各家庭に対して実施された子育てに関するアンケートと同市内における GIS の座標データを用いて、エージェントの行動ルールを If-Then ルールの形式によって決定木を基に生成する。その生成手法において、モデル作成者が一意に決定した前件部を固定して If-Then ルールを生成する手法と、ID3 によって前件部に用いる属性を決定する手法を比較し、各手法におけるアンケートへの正答率などの識別率を比較する。そして結果を基に、ID3 をアンケートからの If-Then ルールの抽出に用いる際に必要な改良点を示す。さらに、前件部に用いる属性値のグルーピングを行うことで、ルール数の減少やアンケートの分析を行う。

2. 子育てに関するアンケート

本研究では、大阪府吹田市内の 6 つの幼稚園のいずれかに属する各家庭に対して実施された子育てに関するアンケートを用いる。アンケートは選択型回答形式であり、その中から任意の問答群をエージェントの属性として If-Then ルールにおける前件部に用いる。If-Then ルールの生成に用いる属性とその属性値を次項に示す。ロコミとはこれまでに他者からの影響によって施設を選択した経験の有無であり、公共施設と民間施設を区別して調査した。また月収と年齢は同一家庭内の夫と妻を異なる属性として用いている。なお、本アンケートの回答者数は 1492 名である。

交通手段

1. 徒歩
2. 自転車
3. 自家用車
4. 電車・バス

通院時間

1. 5分未満
2. 5分～10分未満
3. 10分～15分未満
4. 15分～20分未満
5. 20分～25分未満
6. 25分～30分未満
7. 30分以上

口コミ（公共施設・民間施設）

1. ある
2. なし

月収（妻・夫）

1. 収入なし
2. 5万円未満
3. 10万円未満
4. 10万円～20万円未満
5. 20万円～30万円未満
6. 30万円～40万円未満
7. 40万円～50万円未満
8. 50万円～60万円未満
9. 60万円～70万円未満
10. 70万円～80万円未満
11. 80万円以上

年齢（妻・夫）

1. 20歳代
2. 30歳代
3. 40歳代
4. 50歳代

幼稚園

6つの幼稚園

郵便番号

67の郵便番号

3. 地理情報システム（GIS）

社会シミュレーションはその設計がより現実的であるほど、有効な知見や考察を得ることが出来る。現実的なシミュレーションでの結果から得られた知見や考察を現実社会に還元することで、信頼性の高い政策提案が打ち出される。GISとは現実社会を地図上に抽象化し、多様な属性情報を一元管理することで、地域的な事象を収集、管理、分析、可視化できるシステムである。GISにおける地図はデジタル化された数値地図を基に構成され、データベース化された地理情報と対応している。これらのことから近年、シミュレーション上での仮想環境の構築にGISのデータがしばしば用いられている。GISのデータが用いられた研究として、交通分析による交通モデルの構築などがある[5]。さらに、地震や津波などの自然災害においてその被害範囲や避難経路の予測などの研究がなされている[6,7]。一方、MAS内にGISのデータをどのように組み込むかは課題がある。Goncalvesら[8]はGISとMASの融合システムについて提案している。このシステムでは、GISのデータをあらかじめ取り込み、それを用いてMASを行っている。このことからGISのデータを動的にMASに取り込むことが課題となっている。

本研究では、GISにおける大阪府吹田市の町丁目と医療機関の座標データをあらかじめ取り込む。それらのデータを用いて、エージェントの居住地と医療機関の間の距離計

算を行うことで日常生活では正確に計測することが困難な属性の算出を行った。なお、エージェントは居住地として回答した郵便番号の地区の中心座標に分布することとする。これは吹田市全域で約 35 万人の人口と比較して、本研究で用いるアンケートの回答者が 1492 名であることから、空間におけるエージェントの過分散を避けるためと、アンケートの実施において個人情報保護法などにより、各エージェントの詳細な住所情報は入手や扱いが困難なためである。吹田市内における各郵便番号の中心座標は GIS を用いて算出することができる。

4. 決定木を用いた行動ルールの同定

本研究では、If-Then ルールの生成において決定木を用いる。ここでは、最上位ノードから各節点を経て、最終ノードを結んだものが 1 つの If-Then ルールとなる。したがって、各節点は 2. で示した属性値のいずれかとなる。生成される If-Then ルールは式(1)のように表すことができる。

$$\text{Rule } R: \text{If } A_1 \text{ is } V_{1i} \text{ and } A_2 \text{ is } V_{2i} \text{ and } A_3 \text{ is } V_{3i} \text{ then Class } C \quad (1)$$

ここで R はルールのラベル、 A は属性、 V_{ij} は属性 A に含まれる属性値で、アンケートにおける各エージェントの回答により値が定められる。また C はルールの結論部である。各 If-Then ルールの結論部は医療機関までの距離区分とし、400m 未満、400m から 800m 未満、800m から 2000m 未満、2000m 以上の 4 クラスとする。各ルールの最終ノードにおいて、前件部に使用された属性に基づくエージェントを結論部の 4 クラスに分類し、最大エージェント数を持つ距離区分が結論部となる。この最大エージェント数を持つ結論部クラスが複数ある場合は、結論部を保留とする。識別率は各エージェントが持つ属性に基づいた If-Then ルールの結論部における距離区分と、実際にアンケートで回答している医療機関までの距離区分が同じなら正答、異なる場合は誤答、保留の場合はそのエージェントを保留として計算する。

本稿では、1492 名分のアンケート結果から、2. で示した属性データを全て持つ 497 名のサンプルを用いた実験と 1492 名のサンプル全てを用いた実験を最初に行う。これは、データのエラーやノイズをどのように扱うかが非常に難しい問題となっているためである。アンケートデータにはエラーやノイズと呼ばれるものが多く見られる。これらは、データにおける欠損値や信頼性を意味するものである。郵便番号については、吹田市全域の 95 種類のうち、アンケートから 67 種類の郵便番号が抽出できた。しかし、郵便番号を属性として If-Then ルールの生成に用いる場合にルール数が大幅に増加する。よって郵便番号を除いた 9 属性での実験も行うことで、ルール数の制限を行うとともに医療機関選択における居住区域が与える影響力を調査する。決定木の構築において、各節点を属性値とする場合生成されるルール数は指数乗に増加する。ここでは、ルール数を制限するため、決定木は最大 3 階層に限定して構築する。さらに、月収(夫・妻)の属性値

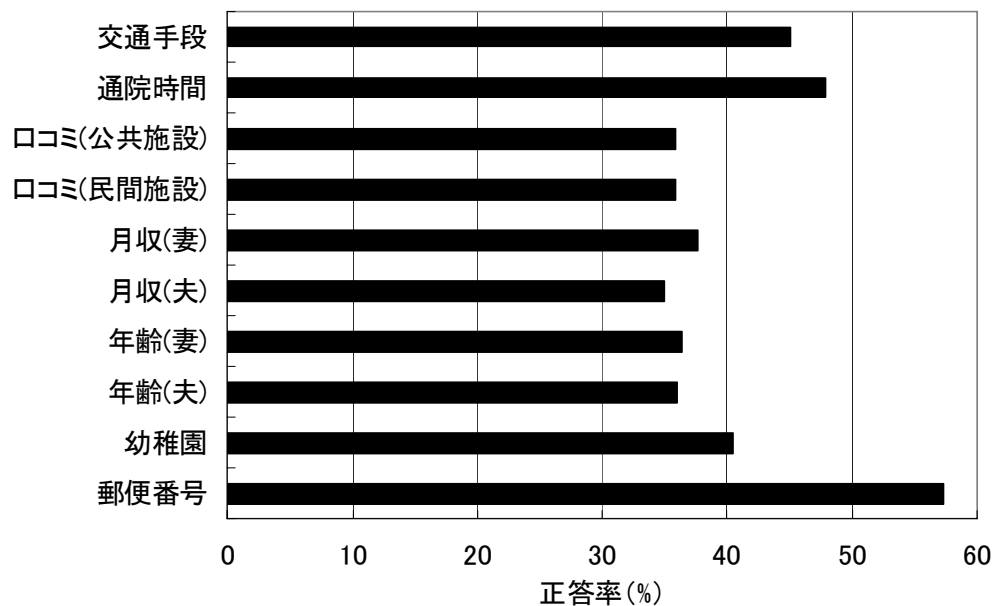


図1 各属性のみでサンプルを分類した識別率

を、以下のように11個から4個に置き換える。

1. 収入なし～10万円未満
2. 10万円～30万円未満
3. 30万円～50万円未満
4. 50万円以上

4.1 固定属性選択

If-Then ルールを生成する上で、条件部に用いる属性を決定するには様々な手法が提案されている。ここではモデル作成者が一意に属性を決定し、識別率の計算対象となる全てのエージェントに対して、決定された属性を条件部に用いて決定木を構築する。本稿ではルール数の増加を制限するため、最大3階層としているので、前件部に用いられる属性の組み合わせは郵便番号を除く9属性での実験において最大 ${}_9C_3 = 84$ 通りあり、全属性を用いる10属性での実験では最大 ${}_{10}C_3 = 120$ 通りある。

ここで、本研究に用いるアンケートにおいて各属性が与える影響を把握するために、1属性のみを用いた正答率を図1に示す。さらに最大3階層として算出した識別率において正答率が最大と最小の結果を表1と表2に示す。

表1と表2から、全属性値を持たないエージェントを含めた実験結果に対して、全属性値を持つエージェントのみの実験結果はいずれも保留率が高くなっている。これは、決定木の最終ノードにおいてエージェント数が極端に減少し、各距離区分のエージェント数が同数になり、保留の結論部を持つIf-Thenルールが生成されやすくなることを示している。さらに、9属性に対して10属性の実験は保留率が上昇する。これも上記と

表 1 全属性値を持たないサンプルを含めた固定属性選択

| | 9 属性 | | 10 属性 | |
|-------|-------|-------|-------|-------|
| | 最大 | 最小 | 最大 | 最小 |
| サンプル数 | 761 個 | 683 個 | 740 個 | 694 個 |
| 正答率 | 56.8% | 32.8% | 81.8% | 59.1% |
| 誤答率 | 37.2% | 60.0% | 8.2% | 24.1% |
| 保留率 | 6.0% | 7.2% | 10.0% | 16.8% |
| ルール数 | 99 個 | 13 個 | 528 個 | 181 個 |

表 2 全属性値を持つサンプルに対する固定属性選択

| | 9 属性 | | 10 属性 | |
|-------|-------|-------|-------|-------|
| | 最大 | 最小 | 最大 | 最小 |
| サンプル数 | 497 個 | 497 個 | 497 個 | 497 個 |
| 正答率 | 56.7% | 31.6% | 80.9% | 59.8% |
| 誤答率 | 36.2% | 54.1% | 5.2% | 18.9% |
| 保留率 | 7.1% | 14.3% | 13.9% | 21.3% |
| ルール数 | 73 個 | 11 個 | 328 個 | 149 個 |

同じ理由であり、郵便番号は 67 種類あるのでノード数が非常に増え、各ノードのエージェント数が極端に減少することが原因と考えられる。しかし、郵便番号を If-Then ルールの生成に使用することで正答率が大幅に上昇している。これは医療機関選択において居住区域が重要な影響力を持つと考えることができる。ここで課題となるのは最終ノードにおいてエージェント数が極端に減少することと、正答率に大きな影響を与える郵便番号の使用方法である。さらに、モデル作成者が一意に属性を決定するので、属性の全ての組み合わせから最適な決定木を選択しなければならない。

4.2 ID3

固定属性選択で生じた問題を解決するために ID3 を用いて決定木を構築する。ID3 は Quinlan[1]によって提案された獲得情報量 (Information gain) を用いて属性の選択を行いながら決定木を構築する手法である。その基本原理は獲得情報量の期待値の最大化である。情報量の概念の導入により、決定木構築に要する分割回数を少なくすることができる。ID3 のアルゴリズムを以下に示す。

サンプル集合を D とする。 D に属する各データはそれぞれ C_1, C_2, \dots, C_n の n 個のクラスのいずれかに属している。各データは A_1, A_2, \dots, A_l で表される l 個の属性を持つ。属性 A_i に対して、 $V_{i1}, V_{i2}, \dots, V_{ij}$ は属性値の集合とする。

獲得情報量 $Gain(A_i, D)$ は次のように求められる。

$$Gain(A_i, D) = I(D) - E(A_i, D) \quad (2)$$

$$I(D) = -\sum_{k=1}^n (p_k \cdot \log_{10} p_k) \quad (3)$$

$$E(A_i, D) = \sum_{j=1}^m p_{ij} \cdot I(D_{ij}) \quad (4)$$

$$p_k = \frac{|D_{C_k}|}{|D|}, \quad p_{ij} = \frac{|D_{F_{ij}}|}{|D_a|} \quad (5)$$

- $Gain(A_i, D)$: 属性 A_i で分類後の相互情報量
 $I(D)$: データ集合 D の情報量
 $E(A_i, D)$: 属性 A_i で分類後の情報量の期待値
 D_{ij} : データ集合 D を V_{ij} に基づいて分割したもの
 P_k : D に属するデータに対して、クラス C_k が出現する確率
 P_{ij} : 属性 A_i に対して、属性値 F_{ij} が出現する確率

上記の式(2)から(5)によって求められる獲得情報量の期待値を最大にする属性を各節点として決定木を構築する。ID3 を用いて算出した識別率を表 3 に示す。さらに、生成された決定木での各階層における属性の使用回数を表 4 に示す。

表 3 において 4.1 での固定属性の結果から、全属性値を持つサンプルのみの結果を示している。表 2 と表 3 の比較から、If-Then ルールの生成におけるアンケートからの属性抽出と決定木の構築に ID3 を用いることで、ルール数を大幅に減少させることができた。これは、固定属性では一意に決定した各属性における全ての属性値が決定木のノードとなることで、ノード数が指数乗に増える。これに対して ID3 では上位ノードの期待値に対して、下位ノードの期待値が小さい場合に、下位ノードへ遷移しないためである。さらに ID3 では属性を固定していないことから、各属性の属性値において、それぞれ該当しているエージェントの結論部クラスへの分布を情報量の概念を基に解析して、柔軟に属性および属性値を決定しているためである。本研究では、結論部クラスへのエージェント分類において最大エージェント数を持つクラスが複数存在する場合に、その If-Then ルールの結論部を保留としている。固定属性では、最終ノードにおけるエージェント数の減少化に伴い、最大エージェント数を持つクラスが複数生じる確率が高くなり、保留率が上昇する。しかし、ID3 による属性選択は下位ノードの探索において、獲得情報量の期待値が高い属性、つまり結論部クラスにおけるエージェントの分布により偏りのある属性を選択するため、保留率を減少させることができた。

表 3 全属性値を持つサンプルに対する ID3

| 正答率 | 誤答率 | 保留率 | ルール数 | サンプル数 |
|-------|-------|-------|-------|-------|
| 66.4% | 17.3% | 16.3% | 151 個 | 497 個 |

表 4 生成された決定木での各階層における属性使用回数

| | 第 1 階層 | 第 2 階層 | 第 3 階層 |
|-----------|--------|--------|--------|
| 交通手段 | 1 | | |
| 通院時間 | | 2 | 25 |
| ロコミ(公共施設) | | 1 | 10 |
| ロコミ(民間施設) | | | 4 |
| 月収(妻) | | | 1 |
| 月収(夫) | | | 6 |
| 年齢(妻) | | | |
| 年齢(夫) | | | |
| 幼稚園 | | | 2 |
| 郵便番号 | | 1 | 3 |

しかし、表 3 の結果を算出する決定木では表 4 で示している通り、第 1 選択属性に交通手段が選択されている。図 1 での各属性の比較において、交通手段の正答率に対して通院時間と郵便番号の正答率の方が大きい。これは、ID3 では結論部である 4 クラスに対して獲得情報量を計算するため、最大エージェント数を持つ結論部クラス以外の結論部クラスのエージェント数が均一である場合に獲得情報量の期待値が小さくなってしまふことがある。

表 5 の例を用いて説明する。属性 A_1 , A_2 について結論部クラス C_1 , C_2 , C_3 , C_4 があると仮定し、それぞれのサンプル数が表 5 のようになっている。ここで情報量が小さいほどデータに偏りがあることとなる。この時 C_1 が正答クラスであるとすると、正答率が高いのは属性 A_1 である。しかし ID3 では、サンプルの分布により偏りがある属性 A_2 が選択される。よって ID3 は正答数最大化を目的とする決定木構築を行う際、最大エージェント数を持つクラス以外のクラスにおけるデータの分布も考慮に入れてしまうという課題がある。

表 5 例題

| | C_1 | C_2 | C_3 | C_4 | 情報量 |
|----------|-------|-------|-------|-------|------|
| 属性 A_1 | 60 人 | 15 人 | 15 人 | 10 人 | 0.48 |
| 属性 A_2 | 40 人 | 35 人 | 25 人 | 0 人 | 0.47 |

表 6 例題の結論部 2 クラス化

| | C_1 | $(C_1 + C_2 + C_3)/3$ | 情報量 |
|----------|-------|-----------------------|------|
| 属性 A_1 | 60 人 | 13.33 人 | 0.25 |
| 属性 A_2 | 40 人 | 20 人 | 0.30 |

表 7 ID3 の結論部 2 クラス化

| 正答率 | 誤答率 | 保留率 | ルール数 | サンプル数 |
|-------|-------|------|-------|-------|
| 80.3% | 15.9% | 3.8% | 169 個 | 497 個 |

表 8 生成された決定木での各階層における属性使用回数

| | 第 1 階層 | 第 2 階層 | 第 3 階層 |
|-----------|--------|--------|--------|
| 交通手段 | | 13 | 47 |
| 通院時間 | | 14 | 15 |
| ロコミ(公共施設) | | 6 | 9 |
| ロコミ(民間施設) | | 2 | 7 |
| 月収(妻) | | 9 | 15 |
| 月収(夫) | | 2 | 2 |
| 年齢(妻) | | 1 | 4 |
| 年齢(夫) | | 5 | 1 |
| 幼稚園 | | 6 | 14 |
| 郵便番号 | 1 | | |

4.3 結論部の 2 クラス化

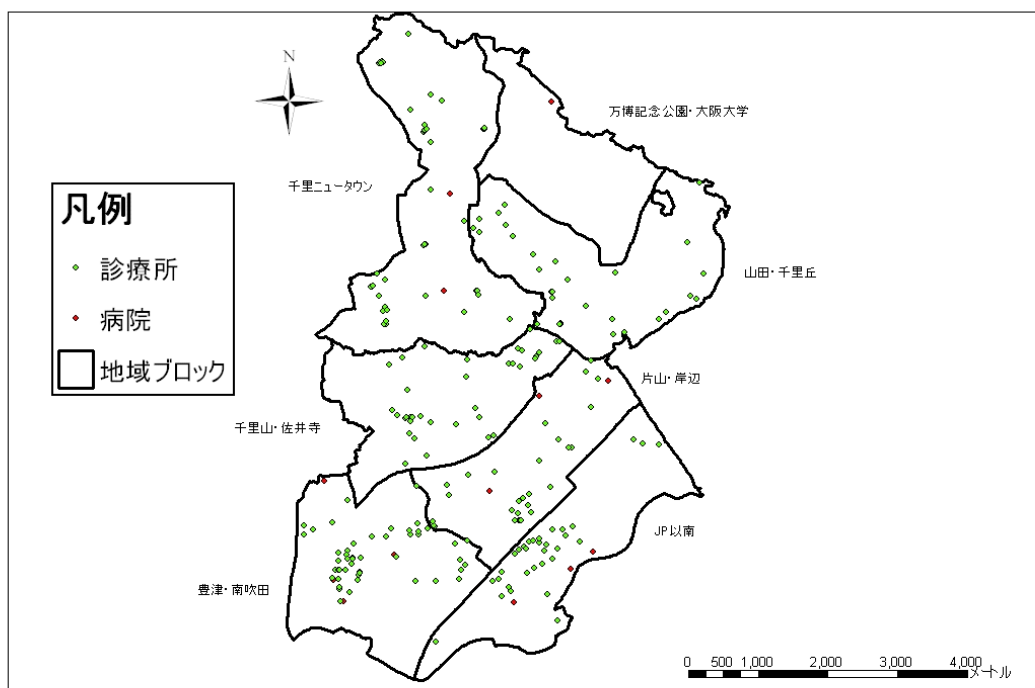
ID3 での課題を解決するために、結論部クラスにおいて最大エージェント数を持つクラスと、それ以外のクラスに分布するエージェント数を平均した 2 クラスにする。この例について表 6 を用いて説明する。表 6 は表 5 での分布を基に結論部 2 クラスへ適用している。最大エージェント数を持つクラス以外の C_2 , C_3 , C_4 を平均化する。これにより最大エージェント数を持つ属性 A_1 が選択される。結論部 2 クラスでの ID3 の結果を表 7 に示す。さらに、生成された決定木での各階層における属性の使用回数を表 8 に示す。表 3 と表 7 の比較から、結論部の 2 クラス化により正答率を大幅に上昇させることができた。これは ID3 における第 1 選択属性に、図 1 において最大正答率を示す郵便番号が選択されたためである。さらに 2 クラス化による効率的な分類により、保留率を減少させることができた。

5. 属性値のグルーピング

5.1 主要7地域

次に我々は医療機関選択への鉄道や幹線道路での地域の隔たりの影響力を調査する。大阪府吹田市は鉄道や幹線道路を基に、大きく7地域に分類することができる。これを図2に示す。4.1で前件部の属性に郵便番号を用いる有効性を示した。しかし郵便番号の属性値が多数であるため、ルール数が多くなる傾向が見られる。したがって郵便番号を主要7地域に分類することで、医療機関選択への鉄道や幹線道路の影響力を調査するとともにルール数の減少を目指す。各エージェントの郵便番号を主要7地域に分類してID3を用いた結果を表9に示す。

表7と表9の比較から正答率が大幅に減少していることがわかる。これは本研究で用いているアンケートにおいて、鉄道や幹線道路での隔たりによる7地域への分類が各地域の特色を薄れさせていると考えられる。したがって有効なルール数の減少ではない。



地域ブロック図

図2 大阪府吹田市における主要7地域

表9 郵便番号を主要7地域へグルーピング

| 正答率 | 誤答率 | 保留率 | ルール数 | サンプル数 |
|-------|-------|------|------|-------|
| 55.7% | 40.0% | 4.3% | 44 個 | 488 個 |

表 10 各属性値の特徴に基づいたグルーピング

| 正答率 | 誤答率 | 保留率 | ルール数 | サンプル数 |
|-------|-------|------|-------|-------|
| 70.4% | 22.7% | 6.9% | 112 個 | 497 個 |

5.2 属性値の特徴に基づいたグルーピング

5.1 における問題に対して、より有効的な属性値のグルーピングを行う。ここでは図 1 で示した 1 属性で生成された If-Then ルールから、結論部が同じ距離区分で、かつ隣接する属性値を新しい属性値としてグルーピングする。グルーピングする属性は、通院時間、月収（夫・妻）、郵便番号である。これらは、各属性値の示している値が連続しているためグルーピングすることができる。郵便番号は、それぞれの隣接情報をあらかじめ作成しておくことで、これを実現している。この方法によって算出した識別率を表 10 に示す。さらに、グルーピングした新しい各属性値を以下に示す。

表 9 と表 10 の比較から正答率が大幅に上昇していることがわかる。これは各属性値におけるサンプルの分類結果を基に、有効な属性値のグルーピングができたと考えられる。しかし表 7 と表 10 の比較から、ルール数の減少とともに正答率の減少も見られる。これは、属性値のグルーピングによりアンケートにおける詳細な情報が粗化されたためだと考えられる。

通院時間

1. 5 分未満 2. 5 分～15 分未満 3. 15 分以上

月収（妻）

1. 5 万円未満 2. 10 万円未満
 3. 10 万円～20 万円未満 4. 20 万円～30 万円未満
 5. 30 万円～40 万円未満 6. 40 万円～60 万円未満
 7. 60 万円～70 万円未満 8. 80 万円以上

月収（夫）

1. 収入なし 2. 5 万円未満
 3. 5 万円～20 万円未満 4. 20 万円～70 万円未満
 5. 70 万円～80 万円未満 6. 80 万円以上

郵便番号

- 67→30

6. おわりに

本研究では、大阪府吹田市における医療機関選択アンケートと GIS の座標データを基に信頼性を考慮した If-Then ルールの生成を行った。If-Then ルールの生成において本稿では決定木構築手法を用いた。決定木の構築において、モデル設計者が一意に決定した属性を用いる手法に対して、ID3 による柔軟な属性選択の有効性を示した。両手法にお

いて決定木を構築し、それを分析することでアンケートからの属性抽出と、それを基にした If-Then ルールの生成に ID3 を用いた場合の問題点を示した。これは、ID3 の特徴である情報量の概念をアンケートデータに用いる場合に、結論部をどのような形にするかが重要であることを示す。そしてこの問題を解決するために、結論部の 2 クラス化を行った。これにより、正答率の大幅な上昇が見られた。さらにルール数の減少を目指し、属性値のグルーピングを行った。ここでは、郵便番号を主要 7 地域にグルーピングする方法に対して、属性値の特徴に基づいたグルーピングの有効性を示した。そしてこれらの実験から、正答率とルール数の関係がトレードオフであることを確認した。今後は、より少ないルール数での正答率の上昇を目指す。

参考文献

- [1] John Ross Quinlan: Induction of decision trees, *Machine Learning*, Vol.1, pp81-106, 1986.
- [2] 馬野元秀: ID3, 日本ファジィ学会誌, Vol.6, No.3, pp502-504, 1994.
- [3] 入月康晴, 古橋武: ファジィエントロピーに基づくファジィ ID3 の提案, 日本ファジィ学会誌, Vol.14, No.3, pp329-333, 2002.
- [4] Jianbing Huo, Xizhao Wang, Mingzhu Lu, Junfen Chen: Induction of Multi-stage decision tree, *Proc. of IEEE International Conference on Systems, Man, and Cybernetics*, pp.835-839, October 8-11, 2006.
- [5] 秋山孝正, 奥嶋政嗣, 和泉範之: マルチエージェント型ファジィ交通行動モデルの提案, 土木計画学研究・論文集, Vol.24, pp.489-498, 2007.
- [6] 中村克行, 小川進: GIS による東海豪雨水害の被害推定, 日本写真測量学会学術講演会発表論文集, Vol.2001, pp.115-116, 2001.
- [7] 宮島宙, 堀宗朗, 小国健二: 地震避難行動シミュレーションのためのマルチエージェントの開発, 応用力学論文集, Vol.10, pp-535-541, 2007.
- [8] Antonio S. Goncalves, Amanda Rodrigues, Luis Corrieia: Multi-agent simulation within geographic information systems, *SCS-The Society for Modeling and Simulation International (United States)*, May, 2004, 6pages.