2006年度 卒業論文

データマイニングアルゴリズム 「アプリオリ」と「ID3」の比較

研究指導 菊池浩明 教授 東海大学 電子情報学部 情報メディア学科 3ADM1127 倉野 奈央子 3ADM3119 阿久津 忍

目次

第1章 はじめに

第2章 データマイニングアルゴリズム

- 2.1 データマイニングとは
- 2.2 アプリオリ
- 2.3 決定木
 - 2. 2. 1 ID3
 - 2.2.2 枝刈り

第3章 実験評価

- 3.1 使用ツール紹介
 - 3.1.1 Gather
 - 3. 1. 2 ID3E
- 3.2 アンケートデータ
- 3.3 アプリオリシステムの性能
- 3.4 出現属性の比較
 - 3.4.1 重要属性の比較
 - 3.4.2 出現属性の適合率
- 3.5 論理関係の比較
- 3.6 考察

第4章 評価ツールの開発

- 4.1 概要
- 4.2 機能
- 4.3 使用方法

第5章 おわりに

参考文献

謝辞

付録1

付録2

第1章

はじめに

現在、情報技術の発達に伴い大量のデータの中から有益な情報、知識を抽出する技術、データマイニングが社会的に必要とされている。データマイニングでは、アプリオリ[1] と ID3[2] という2つの主要なアルゴリズムがしばしば用いられることが多い。しかし、与えられた目的属性についての論理決定木を作る ID3 に対して、アプリオリでは全ての属性間の組み合わせから相関ルールを抽出するという違いがある。そこで、実際のアンケートデータに適用した出力結果を比較し、出現属性や論理関係には相関があるのか確かめることにした。

本研究では、100 名に対してのアンケートを、代表的な手法である決定木学習アルゴリズム ID3 とアプリオリ分析との比較を行い、その出力結果から出現属性や論理関係から相関があることを確かめる。

第2章

データマイニングアルゴリズム

2.1 データマイニングとは

売店の販売データや電話の通話履歴、クレジットカードの利用履歴など、企業に大量に 蓄積されるデータを解析し、その中に潜む項目間の相関関係やパターンなどを探し出す技 術のことで、以下のような手法がある。

≪代表的なデータマイニングアルゴリズム≫

• 決定木

分析の結果を木構造のモデルで出力する方法。

結果が木構造で出力するため理解しやすく高速に処理出来るが、誤差が発生する可能性がありモデル構築でコンピュータに負荷が掛かる。

代表的なアルゴリズムは ID3。

クラスタ分類

複雑な集合を類似したカテゴリーに分けて、それを理解しようとする方法。 探索的手法で容易にデータ入力が出来るが、正しい距離尺度やウェイトを設定する ことが困難なケースが有り結果の解釈が難しい。

代表的なアルゴリズムは K-means 法。

・マーケットバスケット分析

要素の組み合わせから、それらの関連性を探る方法。

強力な探索的手法で明確に結果を理解できる。データ属性が限定的に扱われ問題の 規模に対して指数的に計算量が増大する。

代表的なアルゴリズムはアプリオリアルゴリズム。

2.2 アプリオリ

アプリオリアルゴリズムとは Rakesh Agrawal らによって提案されたアルゴリズムである。本アルゴリズムは、ユーザが確信度のしきい値(最小確信度:minconf) とサポートのしきい値(最小サポート:minsupp) を与えて、最小確信度以上の確信度と最小サポート以上のサポートをもつ全ての相関ルールを効率よく出力する。データベース中の全トランザクションのうち、アイテム集合 X を含むトランザクションの割合を X のサポートといい、supp(X) と表記する。また、ユーザが与えた最小サポート以上のサポートをもつアイテム集合を頻出アイテム集合(frequent itemset) とよぶ。このとき、相関ルール " $X \Rightarrow Y$ " のサポートと確信度は、それぞれ、

 $supp(X \Rightarrow Y) = supp(X \cap Y)$

 $conf(X \Rightarrow Y) = supp(Y \cap Y) / supp(X)$

と定義されている。したがって、条件を満たす相関ルールを求めるためには、 $supp(X \cap Y)$ $\geq minsupp$, かつ、 $supp(X \cap Y)/supp(X) \geq minconf$ であるようなアイテム集合 Y 、X を生成し、そのサポートを求めればよい。

2.3 決定木

データから知識を獲得する方法として、決定木学習アルゴリズム ID3、C4.5 がある. 決定木を構成するために用いるデータの事を学習データと呼ぶ。ここで学習データ D は、m 個の事例 d1,..., dm からなる。事例はそれぞれに属性値を持つ n 個の属性 a1,..., an とターゲット属性であるクラスからなる。k 個クラスがあったとするとクラスは C1,..., Ck と表され、同じクラスを持つ事例によって k 個の部分集合が作られる。決定木のノードには属性、枝には属性値、葉にはクラスが対応する。

2.3.1 ID3

決定木学習アルゴリズム ID3 は、最も簡潔な木を作ることを目的としている。そのために ID3 は、各属性に対して情報量利得を計算し、その値が最も大きな属性から順にノードとして選択し、決定木を構成していく。次に、これらのアルゴリズムについて本研究で重要な点について述べる。

2.3.2 枝刈り

大規模な学習データから決定木を生成するとき、意味的な内容とは無関係な、非常に複雑な木となることがある。そこで、次のパラメータに基づいて木を簡単化する。

1. M (重み)

分類を進めるために必要な事例数の最小値を意味する。この値を大きくすると、枝刈りの効果が大きくなり、木が小さくなる。

2. CF (信頼度)

決定木の葉の誤り率をどこまで許すか設定できる。小さな値ほど大きな枝刈りが行われるが、誤差が大きくなる。

第3章

実験評価

3.1 使用ツール紹介

今回の実験では、先輩方の作ったツールを使わせていただいたので、ここで紹介する。

3.1.1 Gather

2005 年度卒業生、本多淳司「アプリオリアルゴリズムに基づいたデータマイニングツール"Gather"

アプリオリアルゴリズムに基づいたデータマイニングツール。実験により"Gather"は人間の計算時間の1%も満たない時間で計算することが分かったと、述べられている。

3. 1. 2 ID3E

2004 年度卒業生、並木翼「ユーザビィリティの高い GUI ベースの決定木学習ツール "ID3E"|

ユーザビリティの高い GUI ベースの決定木学習ツールを開発。ユーザビリティ評価実験により、開発ツールの総合実験時間は従来ツールの 85%で済むことが分かったと、述べられている。

3.2 アンケートデータ

本研究で用いたアンケートは、東海大学生 100 人を対象に、2006 年 5 月下旬から 7 月上旬にかけて実施した。アンケートの質問を表 1 に示す。これらの 2 択の質問項目を属性として用いる。

表1:アンケートデータ

	F •				
Q1	あなたの性別は?	男		女	
			63		37
Q2	異性の友達が・・・	多い		少ない	
			39		61
Q3	異性間の友情は・・・	成立す	る	成立しない	
			89		11
Q4	海外に行ったことが・・・	ある		ない	
0.5	+1°1 1°11 1°11	1415	51	1.1.5	49
Qo	さびしがり屋ですか?	はい	67	いいえ	22
06	犬と猫ならどっちが好き?		67	 猫	33
QU	人と油ならとうらかがら:	A	68	1世	32
Q7	自分が好き?	好き	- 00	 嫌い	02
Q 7		71 C	48	Mr.	52
Q8	今ダイエットを・・・	している		していない	
			22		78
Q9	感情を表に出すタイプ?	はい		いいえ	
			49		51
Q10	絶叫マシーンが・・・	好き		嫌い	
			70		30
Q11	付き合った経験が・・・	ある		ない	
			84		16
Q12	自分は遠距離恋愛が出来ると・・・	思う		思わない	
			41		59
Q13	やきもち焼きですか?	はい	0.4	いいえ	00
014	七	シカ・十 フ	64	シャナンナン	36
Q14	相手色に・・・	染まる	44	染まらない	56
015	いつでも自分を1番に考えてほしい?	はい	44	いいえ	30
Q13	いっても日がを「歯に与えてはむい!	10.0	42	0.0.7	58
Q16		ある		ない	
	, , , , , , , , , , , , , , , , , , ,		25	U	75
Q17	一人の時間は大切ですか?	はい		いいえ	
			94		6
Q18	ペア商品を持っている(今いない人は持ちたい?)	持ってる	3	持ってない	1
			42		58
Q19	結婚と恋愛は別だと・・・	思う		思わない	
			66		34
Q20	あなたは束縛をする人ですか?	はい		いいえ	
			25		75

3.3 アプリオリシステムの性能

アプリオリシステムの性能を示すため、Java を用いて相関ルールを抽出するシステム"Gather" を用いて、3.2 節のアンケートデータをアプリオリシステムで分析し、最小サポートや最小確信度についての得られるルール数を調査した。この結果を図 1、図 2 に示す。

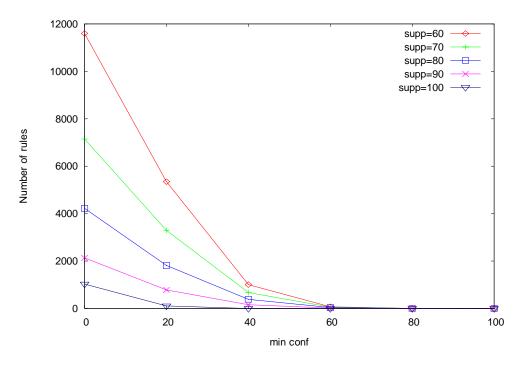


図1: min supp とルール数の関係

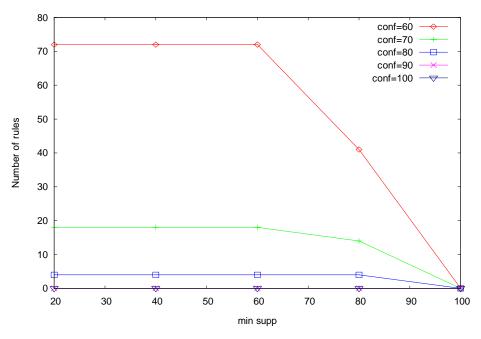


図 2: min conf とルール数

図 1 の結果より、最小確信度の値と得られるルール数には反比例の関係があることが明らかになった。また Conf の値によって、ルール数は大きく変わる。図 2 の結果より、最小サポートの値が一定に達すると、急激にルール数が少なくなることが示された。

3.4 出現属性の比較

ID3 とアプリオリの出力結果による出現属性の信頼性を確かめるため、アンケートデータを ID3 とアプリオリの 2 つのアルゴリズムでデータマイニングし、ID3 とアプリオリそれぞれ、あるターゲット属性について相関のあった属性に関して、2 種類の比較を行う。

1 つ目は、ID3 の出力した決定木で高い位置に出現した属性、アプリオリの出力したルールに何度も出現した属性を重要属性とし、その比較をする。ID3 については枝狩りを行い、出現ノードを減らすことでそれらを重要属性としている。

海外に行った経験があるかどうかというターゲット属性に関してのアプリオリの出力を表 1 に、ID3 の出力を図 3 に示す。また、これらの結果をまとめて表に表したものを表 2 に示す。

表2:海外旅行経験の相関ルール

	公 2 · 1两户广州门门户域、少门口风户					
	出力されたルール					
1	感情を表に出す⇒海外経験ある					
2	自分は遠距離恋愛できない⇒海外経験ある					
3	相手色に染まらない⇒海外経験ある					
4	異性間の友情は成立する△自分は遠距離恋愛できない⇒海外経験ある					
5	ダイエット中でない△相手色に染まらない⇒海外経験ある					
6	一人の時間は大切△相手色に染まらない⇒海外経験ある					
7	束縛しない△相手色に染まらない⇒海外経験ある					
8	感情を表に出さない⇒海外経験ない					
9	感情を表に出さない△一人の時間は大切⇒海外経験ない					

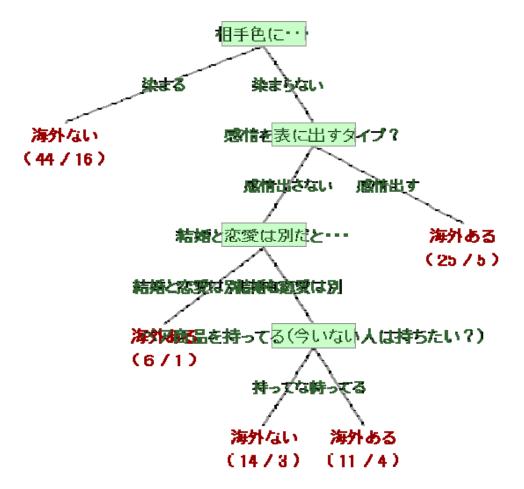


図3:海外旅行経験の決定木 (クラスに分類される事例の数/分類誤りの事例数)

表3:海外旅行経験のアプリオリ、ID3の出現属性

	属性	ターゲット	ID3	アプリ	ーオリ
			m=5,C=9	Supp=0.3,Conf=0.6	
			高さ 4,ノード 4	ルール	√9個
				m=2	m=3
Q1	あなたの性別は?				
Q2	異性の友達が・・・				
$\mathbf{Q}3$	異性間の友情は・・・				1
Q4	海外に行ったことが・・・	0			
Q5	さびしがり屋ですか?				
Q6	犬と猫ならどっちが好き?				
	自分が好き?				
	今ダイエットを・・・				1
Q 9	感情を表に出すタイプ?		2	2	1
Q10	絶叫マシーンが・・・				
Q11	付き合った経験が・・・				
Q12	自分は遠距離恋愛が出来る			1	1
	٤			_	1
	やきもち焼きですか?				
Q14			1	1	3
Q15	いつでも自分をはい番に考				
40-3	えてほしい?				
Q16	人のケータイを黙ってみた				
	ことが・・・				
Q17	一人の時間は大切ですか?				2
Q18	ペア商品を持ってる(今い		4		
	ない人は持ちたい?)		0		
Q19	結婚と恋愛は別だと・・・		3		
Q20	あなたは束縛をする人ですか?				1
	N [→] :				

表 3 は、Q4 をターゲット属性としたときの ID3 とアプリオリの出現属性を表している。 ID3 の列は、属性が出現した属性のノードの高さを表す。アプリオリは、m=2 の列は表 1 のルール 1 のように $X \rightarrow Y$ というように 2 つの属性からなるルール、m=3 の列は表 1 のルール 4 のように $X \land Y \rightarrow Z$ というように 3 つの属性からなるルールの各属性の出現回数を表している。 どちらも出現しなかった属性については空欄となっている。

表 3 を見ると、アプリオリでは Q14 が 4 回、Q9 が 3 回出現している。アプリオリの出力した全ルール数が 9 であるなかで、他の属性が 1,2 回しか出現していないのに対しこれらの属性は数多く出現したことから重要属性であるといえる。また、これらの属性は図 1 の決定木の中でも高い位置に出現している(表 3 では大きな数字となっている)。よって、海外に行った経験があるかどうかという属性に対して、アプリオリと ID3 の重要属性は十分一致したといえる。

また同様に Q9、Q17、Q20 をターゲット属性としたときの ID3 とアプリオリの出現属性を表 4、表 5、表 6 に示す。

表 4: 感情を表に…のアプリオリ、ID3の出現属性

IJ
COC
nf=0.6
個
m=3
1
5
1
1
5
2
2
2
1
1
1
0
3
3
4
3

表 4 の場合においても、Q2 や Q5 はアプリオリで出現回数が多く、ID3 でも出現していることから、重要属性が一致したといえる。アプリオリでは Q19 も出現回数が多いように思うが、m=2 の部分に出現していないことから、minSupport と minConfidence を変えてみれば、アプリオリと ID3 の重要属性の相関関係がもう少し顕著な結果となったのではないかと予想される。

表 5: 一人の時間…のアプリオリ、ID3の出現属性

	双 3 . 八♥フトサラ	ターゲット	スク、ID3の田死海 ID3	アプリ	ノオリ
			m=9,C=24	Supp=0.3	,Conf=0.6
			高さ 4,ノード 4	ルーバ	レ6個
				m=2	m=3
Q1	あなたの性別は?				
$\mathbf{Q}2$	異性の友達が・・・				
$\mathbf{Q}3$	異性間の友情は・・・		1	1	1
Q4	海外に行ったことが・・・				
Q5	さびしがり屋ですか?				
Q6	犬と猫ならどっちが好き?				
$\mathbf{Q7}$	自分が好き?				
$\mathbf{Q8}$	今ダイエットを・・・			1	
Q 9	感情を表に出すタイプ?				
Q10	絶叫マシーンが・・・				
Q11	付き合った経験が・・・		3	1	1
Q12	自分は遠距離恋愛が出来る				
	٤ • • •				
Q13	やきもち焼きですか?		2		
Q14	相手色に・・・				
Q15	いつでも自分をはい番に考				
Q I O	えてほしい?				
Q16	人のケータイを黙ってみた		4	1	
	ことが・・・		1	1	
Q17	一人の時間は大切ですか?	0			
Q18	ペア商品を持ってる(今い				
	ない人は持ちたい?)				
Q19	結婚と恋愛は別だと・・・				
Q20	あなたは束縛をする人です			1	
4.2 0	カュ?			1	

表 5 の結果では、アプリオリと $\mathrm{ID}3$ 両者で出現属性が少なかったが、それらがほぼ一致している。つまり、この場合においても重要属性が一致したといえるだろう。

表 6: 束縛のアプリオリ、ID3の出現属性

	女 0 . 木府		、ID3 少山奶偶性	
		ターゲット	ID3	アプリオリ
			m=9,C=24	Supp=0.3,Conf=0.6
			高さ 5,ノード 5	ルール 25 個
				m=2 m=3
Q1	あなたの性別は?			1
Q2	異性の友達が・・・			1
$\mathbf{Q}3$	異性間の友情は・・・			1
Q4	海外に行ったことが・・・			2
Q5	さびしがり屋ですか?		4	2
Q6	犬と猫ならどっちが好き?		2	1
$\mathbf{Q7}$	自分が好き?			2
Q8	今ダイエットを・・・			1
Q9	感情を表に出すタイプ?		1	2
Q10	絶叫マシーンが・・・			1
Q11	付き合った経験が・・・		3	1
Q12	自分は遠距離恋愛が出来る			2
Q12	٤ • • •			
Q13	やきもち焼きですか?			2
Q14	相手色に・・・			1
Q15	いつでも自分をはい番に考		5	1
QIO	えてほしい?		0	1
Q16	人のケータイを黙ってみた			1
	ことが・・・			1
Q17	一人の時間は大切ですか?			1
Q18	ペア商品を持ってる(今い			1
	ない人は持ちたい?)			
Q19	結婚と恋愛は別だと・・・			1
Q20	あなたは束縛をする人です	\cap		
4.2 0	か?			

表 6 では、アプリオリの出力したルール数が多く、重要属性がわかりにくい結果となってしまった。minSupport と minConfidence の設定を変え、ルール数の調整をして m=3 の場合も求めて比較を行えば重要属性がわかり、他のターゲット属性の場合と同様な結果が得られたのではないかと考えられる。

このように、いくつかのターゲット属性に関しても比較を行ったところ、アプリオリと ID3 の出力するルールの重要な属性は十分一致するといえることが分かった。

2 つ目は、データに対しデータマイニングし、両者の出力に出現した属性を比較、適合率を求める。適合率を求める式を以下に示す。

- ・ アプリオリの適合率=両者に出現した属性数/アプリオリが出力した属性数
- ・ ID3 の適合率=両者に出現した属性数/ID3 が出力した属性数 これらを3つのターゲット属性について求めた結果を表7に示す。

ターゲット属性	アプリオリ	ID3		
Q4	2/7	1/2		
Q7	3/5	3/4		
Q9	3/14	1		
平均	0.37	0.75		

表 7: 出現属性の適合率

アプリオリの平均適合率は 0.37、ID3 の平均適合率は 0.75 となった。全属性数が 20 であることを考慮すると、ターゲット属性 Q7 の例のように出力した 4,5 個の属性のうち半分以上が一致しているというのは高い数字であるといえる。他のターゲット属性ではアプリオリの適合率が低く出てしまったが、これは一般的にアプリオリの方が出力されるルールが多く出るためである。

3.5 論理関係の比較

属性の種類だけではなく、抽出された規則の論理的な矛盾がないかどうかを検討するため次の実験を行った。100 件の集めたデータの中から、Y と N の比率が 50:50 に近い 4 つの属性 (Q4. 海外に行ったことがあるか、Q7. 自分が好きか、Q9. 感情を表に出すタイプか、Q7. 相手色に染まるか)を選び、それらの属性だけからなるデータ D' に対して、アプリオリ、ID3 を適用して一致を調べる。アプリオリの条件は minsupp =0.2、minconf =0.5 で、minconf において高い順に並び替え、その中から答えが海外(決定木のクラス)に関するものだけを抜き出したのが、表 8 である。決定木に対応する葉を最後の列に示しており、括弧は部分的に一致している葉である。

30 · / / / / / / / / / / / / / / / / / /					
ルール	supp	conf	順位	決定木の葉	
感情出す ∧ 染まらない ⇒ 海外ある	0.2	0.8	1位	f	
自分嫌い ∧ 感情出さない ⇒ 海外ない	0.2	0.6897	2位	(e)	
染まる ⇒ 海外ない	0. 28	0.6364	3位	(a,c)	
感情出す ⇒ 海外ある	0.31	0.6327	4位	(f)	
染まらない ⇒ 海外ある	0.35	0.625	5位	(d,f)	
感情出さない ⇒ 海外ない	0.31	0.6078	6位	(e,c)	
自分好き ⇒ 海外ある	0. 29	0.6042	7位	(b)	
自分嫌い ⇒ 海外ない	0.3	0.5769	8位	(a,e)	

表 8: アプリオリから抽出された相関ルール

海外に行ったことが・・・

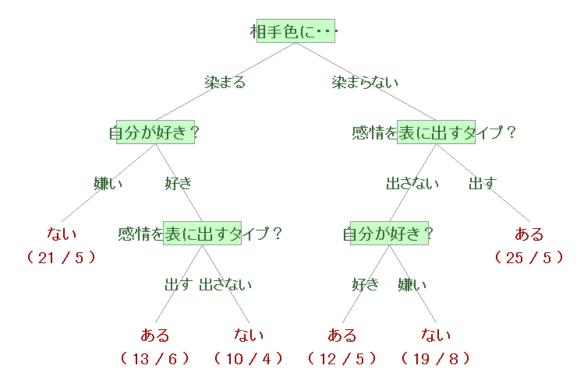


図 4: 「海外旅行経験」の決定木(2)

表 9: 決定木から抽出された相関ルール

	supp	conf
а	0.16	0.76
b	0.07	0.53
С	0.06	0.6
d	0.07	0.58
е	0.11	0.57
f	0.2	0.8

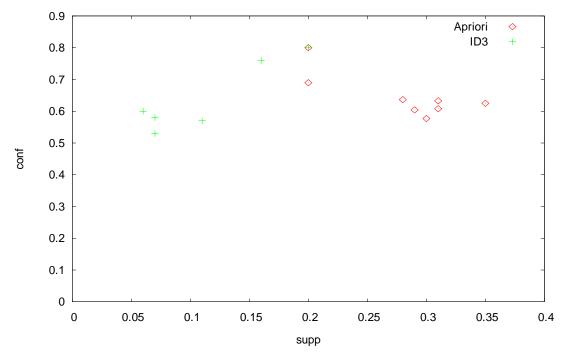


図 5: サポートと確信度についての散布図

次に、図 4 の ID3 の決定木の葉に左から順に番号を振り、supp と conf を表 9 に示した。図 5 は、表 8 から得られた結果から X 軸を supp、Y 軸を conf とする散布図の上に、二種類のルールをプロットした。それらを比較する。図 5 を見て分かるように、アプリオリと ID3 とが一致するもの(表 8 ルール 1 と表 9 の葉 f) があることが分かる。木の枝狩りがうまくいけば、もう少し一致することが予想される。また、ID3 は supp が低く confが高いものが与えられる傾向が高い。アプリオリの条件で supp を下げれば ID3 と一致するより多く生じると考えられる。

2.6 考察

出現属性の比較では、ID3 で高い位置に出現した属性はアプリオリでも出現回数が多かったこと、適合率が高い数値であったことから、両者の出力する属性にはある程度相関があることが分かった。

論理関係の比較では、アプリオリに出力されたルールと ID3 の出力の葉に矛盾がないことから、両者の出力するルールに矛盾が生まれることはないことが分かった。ID3 の葉に一致するアプリオリのルールは少なかったが、決定木から抽出される規則のサポートはアプリオリよりも低いためである。これは、アプリオリの相関ルールがそれぞれ独立なのに対し、決定木は全体で一つの矛盾のない木を構成しなければならないからだと考えられる。

第4章

評価ツールの開発

4.1 概要

両アルゴリズムでの属性の比較にあたり、その一部を自動で行うツールを作成した。

4.2 機能

Java ソース名	機能	入力ファイル (拡張子)	出力ファイ
			ルの拡張子
tocsv.java	Gather の出力ファイルを	Gather の出力ファイル(txt)	csv
	csv ファイルに直す。		
appear.java	出現属性の比較を行う。	属性テーブル(csv)	txt
		tocsv の出力ファイル(csv)、	
		ID3 の出現属性(txt)	
Logic.java	アンケートの回答が半分程	アンケートデータ(csv)	csv
	度に割れた属性部分を抜き		
	出す。		

4.3 使用方法

[出現属性の比較]

・データマイニングする

本ツールは研究室の卒業生のツールを使用し、それらの結果を比較する仕様になっている。まずデータを用意し、Gather、ID3E を使用し、結果(例として図 6、図 3 参照)を保存する。Gatherでは minSupport と minConfidence の設定、ID3E では枝狩りができるが、両者の出現属性の数を本論文を参考に調整してほしい。

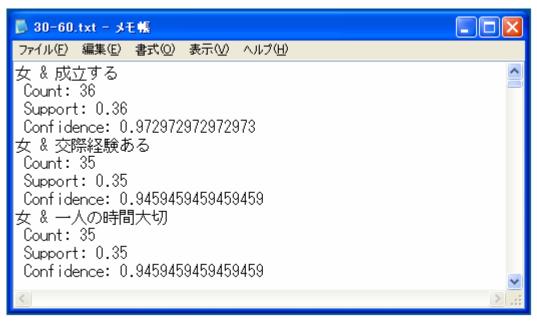


図 6: Gather の出力ファイル例

・アプリオリでの出現属性をまとめる

tocsv.java を実行し、Gather の結果を引数で指定すると、テキストファイルに文章で表示されていた相関ルールが csv ファイルの表にまとめられる。Excel で表示し、相関ルールを「仮定⇒結論」としたとき、結論部の列を昇順(降順でも良い)に並び替え、ターゲット属性が結論となっている相関ルールの仮定部をテキストファイルにコピーアンドペーストする(図 7参照)。

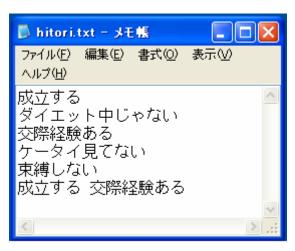


図7:仮定部をコピーしたファイル例

・属性と属性値の表を用意する 属性(質問)と属性値(回答の種類)を表にして保存する(図8参照)。

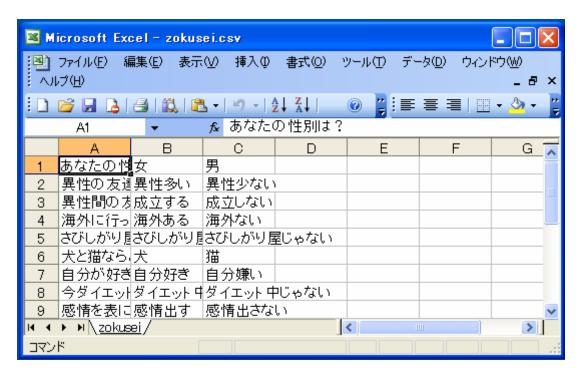


図8:属性テーブル例

・ID3Eでの出現属性をまとめる

ID3E での出現属性をテキストファイルに手動でまとめる。このとき、高い位置のノードから書いていくと、属性の重要さがわかりやすく、結果がわかりやすい。

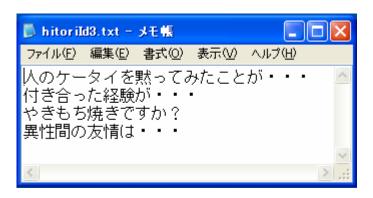


図 9: ID3E 出現属性ファイル例

・比較をする

appear.java を実行し、属性と属性値の表、アプリオリでの出現属性、ID3E での出現属性を引数としてあたえる(図 10 参照)と、各属性のアプリオリでの出現回数、ID3E での出現属性がアプリオリで何回出現しているかが表示され、また同様の内容をアプリオリでの出現属性に「出力結果.txt」が追加されたファイル名で出力する(図 11 参照)。

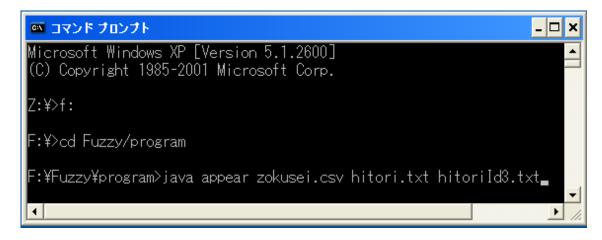


図 10: appear.java 実行画面

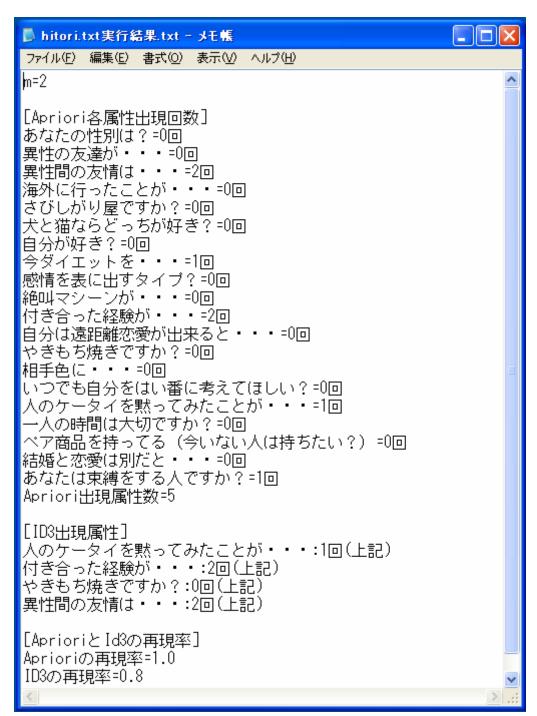


図 11: 出力ファイル例

[論理関係の比較]

論理関係の比較の実験では回答が半分程度($45\sim55$)に分かれている属性を使用している。そこで、Logic.java を実行し入力ファイルを引数で指定する(図 12 参照)と、それらの属性とターゲット属性のみの csv 形式(カンマ区切り)のデータが出力される(図 13 参照)ようになっている。このデータを Gather と ID3E の 2 つのツールにかけ、結果を手動で比較する。ただし、このプログラムは 2 択のアンケートを想定している。

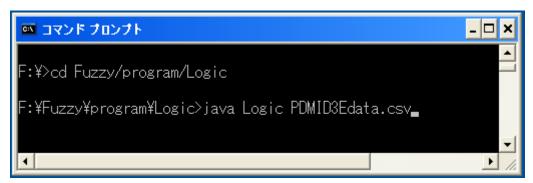


図 12: Logic.java 実行画面

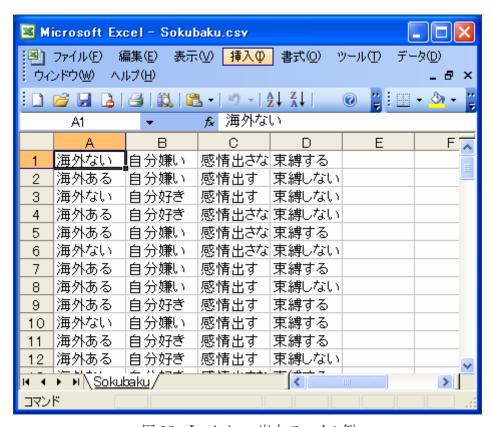


図 13 : Logic.java 出力ファイル例

第5章

おわりに

アプリオリと ID3 の出力結果を比較し、出現属性や論理関係には相関があるのか確かめた。出現する属性に関しては、決定木で上位に出現する属性はアプリオリでも目的属性と相関が高く、平均値でそれぞれ 0.75 と 0.35 の適合率であった。また、抽出された知識の論理的な関係には矛盾はないが、アプリオリで出力された相関ルール 8 つのうち決定木に完全に対応する葉は 1 つしかなかった。これは、アプリオリの相関ルールがそれぞれ独立なのに対して、決定木は全体で一つの矛盾のない木を構成しなければならないところに起因していると考えられる。決定木から抽出される規則のサポートはアプリオリよりもかなり低いことも原因の一つである。これらのことから、アプリオリと ID3 の出力結果に、出現属性や論理関係の相関があると確かめられた。

参考文献

- [1] 本多,アプリオリアルゴリズムに基づいたデータマイニングツール"Gather"の開発,東海大学情報メディア学科2005年度卒業研究論文.
- [2] 松岡, データマイニングによる授業評価アンケートの解析, 日本知能情報ファジィ学会, 気持ちのワークショップ 2003.
- [3] 福田, 森本, 徳山, データサイエンス・シリーズ(3) データマイニング, 共立出版, (2001)
- [4]阿久津, 倉野, 「データマイニングアルゴリズム"アプリオリ"と"ID3"の比較」, 日本知能情報ファジィ学会,曖昧な気持ちに挑むワークショップ,pp. 1-5, 2006

謝辞

本研究を完遂するにあたり、多大なるご指導を受け賜りました東海大学電子情報学部情報メディア学科菊池浩明助教授に心より感謝申し上げます。

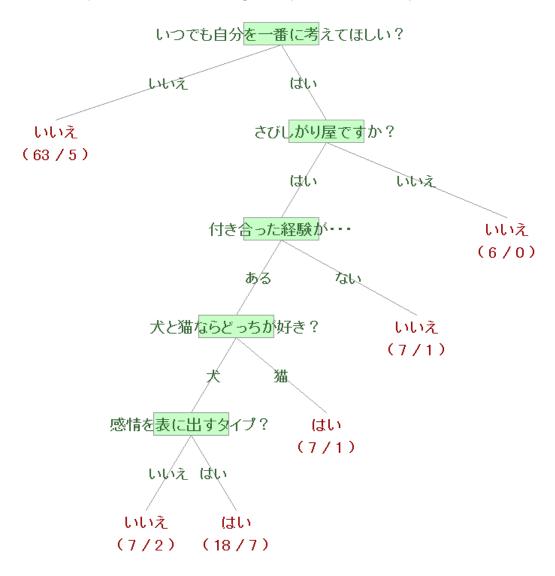
また、多大なるご指導を賜り、本研究を導いて頂きました三橋将氏、ツールを使用させていただいた本多淳司氏、並木翼氏に深くお礼申し上げます。

そして、菊池研究室の同期の皆さんに感謝の意を述べると共に、謝辞とさせて頂きます。

付録1

今回アンケートを取るにあたって、束縛するかどうかをターゲット属性とした。大学生の恋愛感について興味があったためだ。そのID3Eの結果をここに示す。

あなたは束縛をする人ですか?



付録 2

今回のデータマイニングから発見された相関ルールを用いて、Psychology Data Mining という心理テストを夏の課題として作成した。興味のある方は以下のURLよりどうぞ。

http://www.cs.dm.u-tokai.ac.jp/~kitakita/OfficialPage/kadai/index.html