

武藏工業大学	学生会員	高須 光郎
(株) 地崎工業	正会員	須藤 敦史
武藏工業大学	正会員	星谷 勝

1. はじめに

マーケティング分野において、大規模なデータの高度な処理を行うデータマイニング^①の数々の適用事例が発表されている。データマイニングは、膨大なデータの中に存在する隠れた法則や相関関係を自動的に発見する手法である。そして、定性および定量データが混在するデータや属性の相関関係が複雑なデータも扱うことが可能であり、その適用分野を問わないのが特徴である。加えて、土木分野においてもデータマイニングは情報分析手法として対応でき、その有用性も期待される。そこで本研究はデータマイニングの適用を目的とし、その概要や適用例を示すとともに、データマイニングにおいて重要な要素となる条件付きエントロピーの検討を簡単な数値解析を通して行っている。

2. データマイニング

データマイニングの名称は鉱山から鉱脈を探掘することに由来し、93年に米IBMアルマデン研究所のアグラワラ研究員が提唱したものである。従来の解析手法はデータが膨大になると分析が難しくなり、細部に至るまでは人間の目が行き届かず、データが十分に活用できないことがある。しかし、データマイニングは人間では想像もできないような規則性を見つけ出し、データ中に存在する情報を引き出すことが可能である。さらに分析作業を自動化することにより、分析結果の客観性と信頼性を高め、エンドユーザーに高度な意思決定を支援することができる。

適用事例（マーケティング） あるスーパーマーケットにおいて売上実績データをデータマイニングツールで処理し、顧客の購入心理とその規則性について分析した結果、予想以上に多くの既婚男性客がビールと紙おむつと一緒に買っていくパターンを見出した。そこでビールと紙おむつの売場を並べてディスプレイしたところ、売上げの増加につながった例が挙げられる。

データマイニングにおける事象の関係「IF～， THEN …」形式が「事象 a_i が与えられたときの事象 b_j の条件付き確率」に対応している。ここで本研究は情報エントロピーを導入し、条件付き確率と情報エントロピーがデータマイニングの手段としてなりうることを検討する。

3. 条件付きエントロピー

エントロピー^②とは「不確定の度合い」といった極めて抽象的なものを定量的に評価するものであり、ある事象のエントロピーはその確率分布によって一意的に定まる。離散型確率分布 $P(A=a_i)=p_i$ ($i=1,\dots,m$) を有する事象 A のエントロピーは式(1)のように与えられる。

$$H(A) = H(p_1, \dots, p_m) = -\sum_{i=1}^m p_i \log p_i \quad (1)$$

次に、確率分布 $P(B=b_j)=q_j$ ($j=1,\dots,n$) を有する事象 B が起こったときの、 A の条件付きエントロピーは式(2)のようになる。

$$H(A|b_j) = H(P(a_1|b_j), P(a_2|b_j), \dots, P(a_n|b_j)) \quad (2)$$

ここで、 $P(a_i|b_j)$ は b_j が発生したときの a_i の条件付き確率であり、 $H(A|b_j)$ は b_j が発生したときの A の条件付きエントロピーである。これは A と b_j が独立のときは $H(A)$ と一致し、相関がある場合には $H(A)$ に比べ減少する。

KEYWORD データ処理 確率論 条件付き確率 情報エントロピー

連絡先 住所 〒158-0087 世田谷区玉堤1-28-1 TEL 03-3703-3111(ex.3268) FAX 03-5707-2187

4. 数値計算例

図-1から図-3に示す3つの（交通）データ Sample I, II, III を用いて数値解析を行った。この図は運転手の性別 A と速度 B について、男性 215 台と女性 25 台（計 240 台）に分けてヒストグラムで表したものである。

これら3つのデータは、男女合せた全体の総数（240 台）とその分布形は同一である。異なる点は、男女それぞれの分布形である（ただし、データ数は同じである）。Sample I は女性分布がばらついているもの、Sample II は女性分布が寄せ集まっているもの、Sample III は Sample II と同一の女性分布が右に寄っているものである。もちろん男性分布もそれぞれ異なっているが、その特徴についてはここでは省略する。表-1 に男女合せた全体の速度のエントロピー $H(B)$ 、男性であることが確認されたときの速度のエントロピー $H(B|a_1)$ 、女性であることが確認されたときの速度のエントロピー $H(B|a_2)$ 、性別（男女どちらかはわからなくてよい）が確認されたときの速度のエントロピー $H(B|A)$ を示す。 $H(B|A)$ は $H(B|a_1)$ と $H(B|a_2)$ の重み付き平均である。

Sample I ~ Sample III すべてにおいて、 $H(B|A)$ が $H(B)$ に比べ減少している。これより、性別 A と速度 B は独立なものではなく、何らかの相関があることがわかる。またその相関の度合いは、 $H(B)$ と $H(B|A)$ の差より Sample II が最も弱く、以下 Sample I, Sample III の順に強くなる。

次に、 $H(B|A)$ の要素である $H(B|a_1)$ と $H(B|a_2)$ の大小関係について考察する。Sample I ~ Sample III すべてにおいて、 $H(B|a_1) > H(B|a_2)$ となっている。これより、男性に比べて女性のほうが速度との相関が強いことがわかる。

最後に、 $H(B|a_2)$ の変化について考察する。Sample I と Sample II について比較すると Sample II のほうが小さい。これより、女性分布が Sample I のばらついたものより Sample II のまとめたもののほうが $H(B|a_2)$ が小さくなることがわかる。ところが、Sample II と Sample III について比較すると $H(B|a_2)$ に変化が見られない。これより、女性分布の平均値は $H(B|a_2)$ には影響しないことがわかる。つまり、 $H(B|a_2)$ は女性の分布形のみに依存し、平均値には依存しないのである。

5.まとめ

将来、データマイニングは土木分野の中でも多種多様な応用範囲が考えられる。今後は実データに適用する予定である。

<参考文献>

- 1). 情報・通信新語辞典 97 年度版、日経 B P 社、1996 2). 有本卓：確率・情報・エントロピー、森北出版、1992

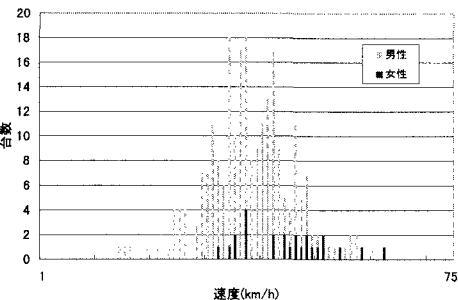


図-1 Sample I

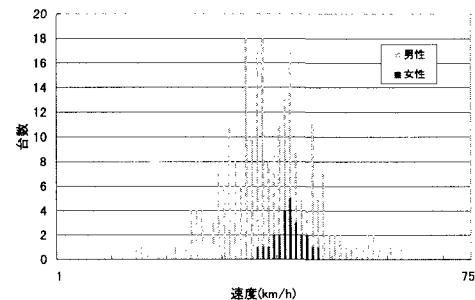


図-2 Sample II

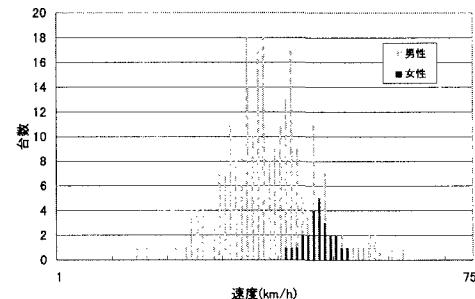


図-3 Sample III

表-1 条件付エントロピーの比較

	$H(B)$	$H(B a_1)$	$H(B a_2)$	$H(B A)$
Sample I	3.36055	3.29113	2.66436	3.22584
Sample II	3.36055	3.37211	2.32154	3.26267
Sample III	3.36055	3.24782	2.32154	3.15133

a_1 : 男性, a_2 : 女性