

# データマイニング

## 第4回授業

担当：菊池

# 相関ルール抽出問題

## 【相関ルール抽出問題】

- ・ 確信度の閾値（最小確信度）
- ・ サポートの閾値（最小サポート）

の2つを入力したとき，それぞれの閾値以上の確信度・サポートを有する相関ルールを全て発見せよ．

条件を満たす相関ルールの抽出問題

・・・下記の2つの部分問題に分解された．

- (1) 頻出アイテム集合を全て見出し，サポートを求める（アプリアリ・アルゴリズム）
- (2) 上記(1)で求めた頻出アイテム集合を使って，最小確信度以上の相関ルールを求める

# 相関ルールの導出

これまでは相関ルール抽出問題のうち，第1段階の頻出アイテム集合の導出の処理について説明した．

引き続き，第2段階の頻出アイテム集合を使った最小確信度以上の相関ルールの導出法を紹介する．

まず，効率の良い導出の為に，以下で最小確信度に関する性質を導く．既習の通り，アイテム集合 $a$ の任意の部分集合 $\tilde{a} \subset a$ について，必ず

$$\text{support}(\tilde{a}) \geq \text{support}(a) \quad (\text{i})$$

が成り立つ．

※記号 $\sim$ は「チルダ」または「チルド」と読みます

# 相関ルールの導出

- $\text{support}(\tilde{a}) \geq \text{support}(a)$  (※) の具体例

$a = \{\text{牛乳}, \text{パン}, \text{チョコレート}\}$

$\tilde{a} = \{\text{牛乳}\}$

とすると,

$$\tilde{a} \subset a$$

が成り立っている.

このとき, 「少なくとも牛乳を買う客」の数のほうが, 「少なくとも牛乳とパンとチョコレートの3つを買う客」より明らかに多い.

よって,

「全部の客に対する, 少なくとも牛乳を買う客の割合」

$\geq$  「全部の客に対する, 少なくとも牛乳とパンとチョコレートの3つを買う客の割合」

が成り立つ. この不等式が(※)式に対応する.

$$\text{conf}(a \Rightarrow (l - a)) = \text{support}(l) / \text{support}(a)$$

ここで,  $a$ を頻出アイテム集合の1つとし,  $a \subset l$ とする.

$$a \Rightarrow (l - a)$$

の確信度 $\text{conf}(a \Rightarrow (l - a))$ は

$$\text{conf}(a \Rightarrow (l - a)) = \text{support}(l) / \text{support}(a) \quad (\text{ii})$$

となる.

<(ii)式の成立の理由>

第3回授業より,  $\text{conf}(X \Rightarrow Y) = \text{support}(X \cup Y) / \text{support}(X)$ だった.  
この式に $X$ として $a$ を,  $Y$ として $(l - a)$ を当てはめて考えたとき,

$$\text{conf}(a \Rightarrow (l - a)) = \text{support}(a \cup (l - a)) / \text{support}(a)$$

ここで,

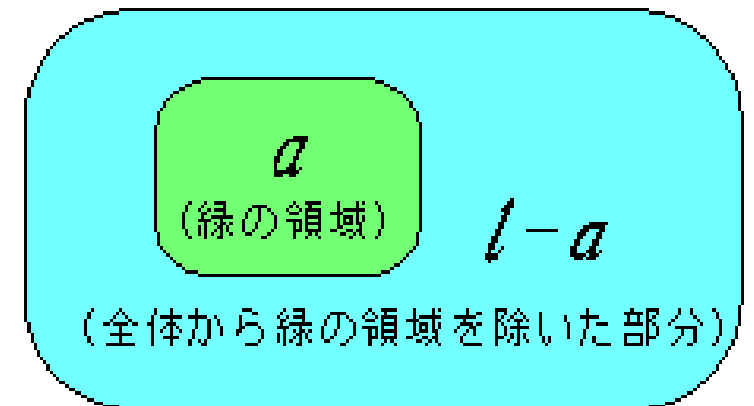
$$a \cup (l - a) = l$$

となるから, (ii)式

$$\text{conf}(a \Rightarrow (l - a)) = \text{support}(l) / \text{support}(a)$$

が成り立つ.

$l$  (全体: 緑の領域と水色の領域の和)



「 $a \Rightarrow (l - a)$ 」の確信度  $\geq$  「 $\tilde{a} \Rightarrow (l - \tilde{a})$ 」の確信度

(ii)式で表された確信度  $\text{conf}(a \Rightarrow (l - a)) = \text{support}(l) / \text{support}(a)$  は,  
(i)より  $\text{support}(\tilde{a}) \geq \text{support}(a)$  なので,

$$\tilde{a} \Rightarrow (l - \tilde{a})$$


の確信度

$$\text{conf}(\tilde{a} \Rightarrow (l - \tilde{a})) = \text{support}(l) / \text{support}(\tilde{a}) \quad (\text{iii})$$

と等しいか、あるいは大きい.

(↑ (ii), (iii)両式の右辺の分母に(i)を適用して考えてみると、明らかに成り立つことがわかる)

すなわち,


$$\text{support}(\tilde{a}) \geq \text{support}(a)$$

「 $a \Rightarrow (l - a)$ 」の確信度  $\geq$  「 $\tilde{a} \Rightarrow (l - \tilde{a})$ 」の確信度  $\cdots \star$

となる.

# $\text{conf}(a \Rightarrow (l - a)) \geq \text{conf}(\tilde{a} \Rightarrow (l - \tilde{a}))$ の具体例

<具体例>

$l = \{\text{牛乳}, \text{パン}, \text{チョコレート}, \text{カップラーメン}\}$

$a = \{\text{牛乳}, \text{パン}\}$

$\tilde{a} = \{\text{牛乳}\}$

とすると,

$$\tilde{a} \subset a \subset l$$

が成り立っている.

このとき  $a \Rightarrow (l - a)$  の確信度とは,  $(l - a) = \{\text{チョコレート}, \text{カップラーメン}\}$  なので,

「少なくとも牛乳とパンを買う客が, どれだけチョコレートとカップラーメンも一緒に買うか」を示す.

いっぽう,  $\tilde{a} \Rightarrow (l - \tilde{a})$  の確信度とは,  $(l - \tilde{a}) = \{\text{パン}, \text{チョコレート}, \text{カップラーメン}\}$  なので,

「少なくとも牛乳を買う客が, どれだけパンとチョコレートとカップラーメンも一緒に買うか」を示す.

上述の黄色の☆で示される関係を適用すると,

「少なくとも牛乳とパンを買う客のうち, チョコレートとカップラーメンも一緒に買う客の割合」のほうが,

「少なくとも牛乳を買う客のうち, パンとチョコレートとカップラーメンも一緒に買う客の割合」より**大きい**ことになる.

# 確信度の性質(1)

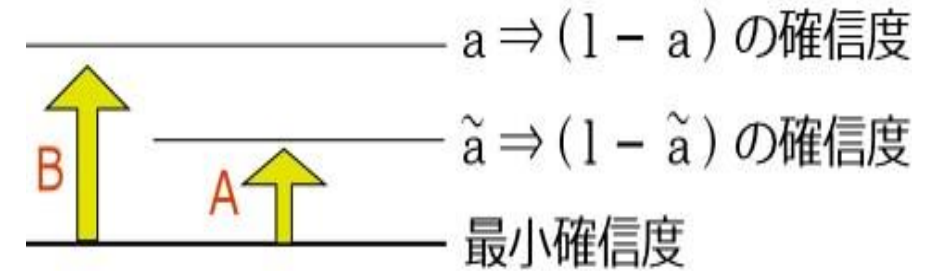
以上より，次の命題が成り立つ．

もし

$\tilde{a} \Rightarrow (1 - \tilde{a})$  の確信度  $\geq$  最小確信度  
ならば，

$a \Rightarrow (1 - a)$  の確信度  $\geq$  最小確信度  
となる．

(※)



「AならばB」が成り立つ

上記は，確信度の低いほうのルールが最小確信度以上であれば，確信度の高いほうのルールも必然的に最小確信度以上になることを示す．



# 確信度の性質(2)

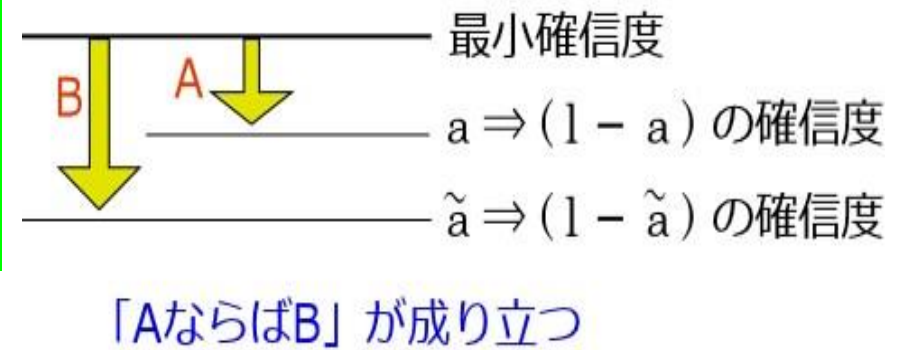
この(※)の命題の対偶をとる(「AならばB」の対偶は「not Bならばnot A」. ある命題が真であれば, その対偶も真になる)と, 次のようになる.

もし

$a \Rightarrow (1 - a)$  の確信度  $\leq$  最小確信度  
ならば,

$\tilde{a} \Rightarrow (1 - \tilde{a})$  の確信度  $\leq$  最小確信度となる.

(※※)



これは確信度の高いほうのルールが最小確信度以下であれば, 確信度の低いほうのルールも最小確信度以下になることを示す.

# 確信度の性質(3)

ここで

$$c = l - \tilde{a}$$

$$\tilde{c} = l - a$$

のようにおく.

このとき  $\tilde{c} \subset c$  が成り立つ(図参照).

上で定義した  $c$ ,  $\tilde{c}$  を用いると,

$$a = l - \tilde{c}$$

$$\tilde{a} = l - c$$

$$l - a = \tilde{c}$$

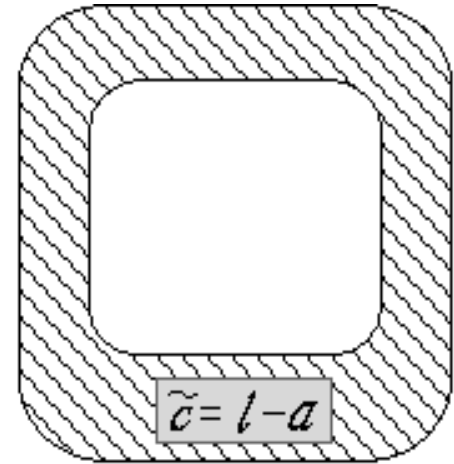
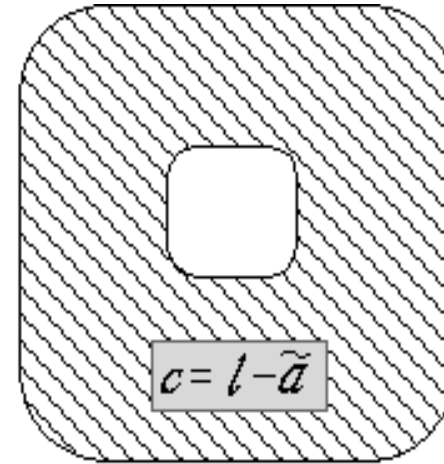
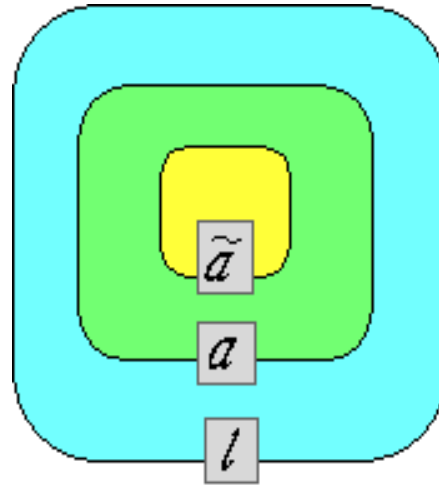
$$l - \tilde{a} = c$$

となるので, 上述の命題 (※※) は次のように表せる.

もし

$(l - \tilde{c}) \Rightarrow \tilde{c}$  の確信度  $\leq$  最小確信度  
ならば,

$(l - c) \Rightarrow c$  の確信度  $\leq$  最小確信度  
となる.



いずれも中央部をくり抜いた斜線部分

前のスライドの(※※)の命題:

もし

$a \Rightarrow (l - a)$  の確信度  $\leq$  最小確信度

ならば,

$\tilde{a} \Rightarrow (l - \tilde{a})$  の確信度  $\leq$  最小確信度となる.

(※※※)

# 相関ルールの確信度に関する重要な性質

以上より、**あるルールが最小確信度未満であれば、そのルールの帰結部の集合を包含するような集合を帰結部に持つ一連のルールも必ず最小確信度未満になる**ので、最小確信度以上のルールの導出の際にはそれらを却下できる。

< 具体例 >

上の(※※※)の命題にて、 $l = \{A, B, C, D\}$ ,  $c = \{C, D\}$ ,  $\tilde{c} = \{D\} \subset c$ とすると

$\{A, B, C\} \Rightarrow \{D\}$  の確信度が最小確信度未満

であるとき、

$\{A, B\} \Rightarrow \{C, D\}$  の確信度は最小確信度未満

となる。

上記の性質を利用して最小確信度(*minconf*)以上の確信度を持つ相関ルールを効率良く導出するアルゴリズムを以下に紹介する

# 最小確信度( $minconf$ )以上の確信度を持つ相関ルール導出のアルゴリズム(1)

## <処理の流れ>

頻出アイテム集合  $I$  について, まず要素数1の結論部分を持つ相関ルールを作成する.  
次に, 確信度が最小確信度より小さくなかった結論部分の集合に対して,  
AprioriGen()関数を適用して, 要素数2の結論部分を生成する.  
以上を繰り返して, 徐々に大きな結論部分を有する相関ルールを求めてゆく.

0) Algorithm Generate-Rules() {

1) foreach 頻出アイテム集合  $I_k$  (要素数  $k > 2$ ) {

2)  $H_1 := \{ \{ h \in I_k \} \mid \text{conf}(I_k - \{h\} \Rightarrow \{h\}) \geq minconf \};$

//  $\uparrow H_m$  (ここでは  $m=1$ ) は最小確信度以上の相関ルールを作れる要素数  $m$  の結論部の集合

//  $h$  は  $I_k$  に含まれる1個のアイテム.

3) call AP-GenRule( $I_k, H_1$ );

4) }

5) }

## 最小確信度( $minconf$ )以上の確信度を持つ相関ルール導出のアルゴリズム(2)

```
6) Procedure AP-GenRule( Itemset  $I_k$ , 結論部の集合  $H_m$ ) {
7)   if ( $k > m + 1$ ) {
8)      $H_{m+1} = \text{AprioriGen}(H_m)$ ;
      // ↑  $H_m$ から要素数の1つ多い $H_{m+1}$ を作り出す. 前回授業で既出
9)     foreach  $h_{m+1} \in H_{m+1}$  {
      // ↑  $H_{m+1}$ に含まれる, 要素数 $m+1$ のあらゆる部分集合 $h_{m+1}$ について以下を実行
10)       $conf = \text{support}(I_k) / \text{support}(I_k - h_{m+1})$ ;
      // ↑ 相関ルール  $(I_k - h_{m+1}) \Rightarrow h_{m+1}$  の確信度を計算
11)      if ( $conf \geq minconf$ ) // 確信度が最小確信度以上か?
12)        output  $(I_k - h_{m+1}) \Rightarrow h_{m+1}$ ;
      // ↑ 相関ルール  $(I_k - h_{m+1}) \Rightarrow h_{m+1}$  は条件を満たすので, このルールを出力
13)      else
14)         $H_{m+1} = H_{m+1} - \{h_{m+1}\}$ ; // 最小確信度以上にはならなかったから,  $h_{m+1}$ を候補から除く
15)    }
16)    AP-GenRule( $I_k$ ,  $H_{m+1}$ ); // 再帰呼び出しで, さらに要素数の多い結論部分について吟味.
17)  }
18) }
```

# 実行例

第3回授業で扱った下記のデータベースより頻出アイテム集合が求まっているとする。  
上記データベースより求まる頻出アイテム集合(最小サポート50%)・・・ $\{\{A\}, \{B\}, \{C\}, \{E\}, \{A, C\}, \{B, C\}, \{B, E\}, \{C, E\}, \{B, C, E\}\}$   
これら頻出アイテム集合に対して“Generate-Rules()”を適用する。  
要素数2より大きい各頻出アイテム集合 $l_k$ に対して、 $H_1$ を生成しつつ、手続き“AP-GenRule”を呼び出す。  
ここでは $l_3$ の(唯一の)要素である $\{B, C, E\}$ に対して上記ルールを適用してみる。

$H_1 = \{\{B\}, \{C\}, \{E\}\}$ .

<確認>

- ・  $h_1 = \{B\}$  に対して： $\text{conf}(\{C, E\} \Rightarrow \{B\}) = \text{support}\{B, C, E\} / \text{support}\{C, E\} = (2/4)/(2/4) = 1 > 0.5$
- ・  $h_1 = \{C\}$  に対して： $\text{conf}(\{B, E\} \Rightarrow \{C\}) = \text{support}\{B, C, E\} / \text{support}\{B, E\} = (2/4)/(3/4) = 2/3 > 0.5$
- ・  $h_1 = \{E\}$  に対して： $\text{conf}(\{B, C\} \Rightarrow \{E\}) = \text{support}\{B, C, E\} / \text{support}\{B, C\} = (2/4)/(2/4) = 1 > 0.5$

よって、 $\{B\}, \{C\}, \{E\}$ は $H_1$ の要素になれる。

上記 $H_1$ より、AP-GenRuleを呼び、

$H_2 = \{\{B, C\}, \{B, E\}, \{C, E\}\}$

が得られる。ここで、最小確信度は0.5とする。

以下、 $H_2$ の各要素に対してルール $(l_k - h_{m+1}) \Rightarrow h_{m+1}$ を吟味する。

・ $h_2 = \{B, C\}$ に対して： $\text{conf}(\{E\} \Rightarrow \{B, C\}) = \text{support}\{B, C, E\} / \text{support}\{E\}$ $= (2/4) / (3/4) = 2/3 > 0.5$ $\therefore \{E\} \Rightarrow \{B, C\}$ は求めるべきルールに含められる。(i)	・ $h_2 = \{B, E\}$ に対して： $\text{conf}(\{C\} \Rightarrow \{B, E\}) = \text{support}\{B, C, E\} / \text{support}\{C\}$ $= (2/4) / (3/4) = 2/3 > 0.5$ $\therefore \{C\} \Rightarrow \{B, E\}$ は求めるべきルールに含められる。(ii)	・ $h_2 = \{C, E\}$ に対して： $\text{conf}(\{B\} \Rightarrow \{C, E\}) = \text{support}\{B, C, E\} / \text{support}\{B\}$ $= (2/4) / (3/4) = 2/3 > 0.5$ $\therefore \{B\} \Rightarrow \{C, E\}$ は求めるべきルールに含められる。(iii)
---	--	---

再度、AP-GenRuleを呼ぶと、 $k=3$ に対して $m+1$ も3になり、となり、条件の $k > m+1$ を満たさなくなるので、ここで終了。

よって、(i), (ii), (iii)が求まるルールとなる。

データベースD	
T I D	アイテム
0001	A, C, D
0002	B, C, E
0003	A, B, C, E
0004	B, E