# Statistical Decision making of Association Threshold in Association Rule Data Mining

*Hee Chang Park[1], Geum Min Song[2]*

## Abstract

One of the well-studied problems in data mining is the search for association rules. In this paper we consider the statistical decision making of association threshold in association rule. A chi-squared statistic is used to find minimum association threshold. We can calculate the range of the value that two item sets are occurred simultaneously, and can find the minimum confidence threshold values.

Keywords :                ,              ,            ,

## 1.

(data mining)                          (mine)

.    ,

.

.                  (association rule)

,               (decision tree),          (neural network),

(clustering),                (genetic algorithm),                        (bayesian network),        -

(memory-based reasoning)                          .

(support),          (confidence),          (lift)

.

,

.

Agrawal    (1993)                          , Agrawal    (1994)

,

[1]Professor, Department of Statistics, Changwon National University, Changwon, Kyungnam 641-773, Korea
E-mail : hcpark @sarim.changwon.ac.kr
[2]Graduate Student, Department of Statistics, Changwon National University, Changwon, Kyungnam 641-773, Korea.

Apriori, AprioriTid                                    . Park        (1995)

partitioning                                    ,  Toivonen(1996)

sampling                                    .            Cheung        (1996)

FUP(fast update)

,  Sergey       (1997)

DIC(dynamic itemset counting)

. Liu        (1999)

DHP(direct hashing and pruning)                                    . Saygin        (2002)

.

Fast

.                                    Silverstein        (1997)

(1)                                    ,

(0)                                    .

(minimum support),

(minimum confidence)                                    ,

,

.

,                                    .

2                                    ,        3

,                                    ,

,

.        4

,

.                        5

.

2.

,

.

.

(transaction) :                                          ,                                .

 (item set) :                                          ,                              .

        (candidate item set) :

                                                              .

        (frequent item set) :

                                                              .

   k-                     $\mathfrak{I} = \{i_1, i_2, \ldots, i_k\}$          $T \subseteq \mathfrak{I}$                              ,                    $T$
TID                                          .          $A \subseteq T$                              $T$      $A$

            .

            $R : X \Rightarrow Y$                              .

            :  $Sup(X \Rightarrow Y) = P(X \bigcap Y)$

            :  $Conf(X \Rightarrow Y) = P(Y|X) = \dfrac{P(X \bigcap Y)}{P(X)}$

            :  $Lift(X \Rightarrow Y) = \dfrac{P(Y|X)}{P(Y)} = \dfrac{P(X \bigcap Y)}{P(X)P(Y)}$

   $X \subset T, Y \subset T, \text{and} X \bigcap Y = \varnothing$        .

            $R : X \Rightarrow Y [support = s\%, \ confidence = c\%]$      X      Y

   s%        , X                              X    Y                                      c%                    .

                                          ,

(minimum support threshold : min_sup)                          (minimum confidence threshold :
min_conf)                                          (item-sets)

      .

            ,                                                                            .

                                                                      ,

            (Lift)                              .

                  Apriori            k-                              k-1

   (join step)                    (candidate k-itemsets)                .

                  ,

(prune step)

.

3.

3.1

(association)

(test of independence)                .                                    2× 2

(contingency table)                        .

| <    1>                     |     |           | 2× 2                              |                       |
| --------------------------- | --- | --------- | --------------------------------- | --------------------- |
|                             |     | Y         |                                   |                       |
|                             |     | 1         | 0                                 |                       |
| X                           | 1   | $a$       | $x_1 - a$                         | $x_1$                 |
|                             | 0   | $y_1 - a$ | $t - (x_1 + y_1) + a$, $x_0 = t - x_1$ |                  |
|                             |     | $y_1$     | $y_0 = t - y_1$                   | $t$                   |

$t$ :

$x_1$ :              X

$y_1$ :              Y

$a$ :              X  Y

$x_0 = t - x_1$ :              X

$y_0 = t - y_1$ :              Y

.                    (3.1)                            .

$$
\left.
\begin{array}{l}
0 \le a \le t \\
0 \le x_1 - a \le t \\
0 \le y_1 - a \le t \\
0 \le t - (x_1 + y_1) + a \le t \\
0 \le a \le x_1 \\
0 \le a \le y_1
\end{array}
\right\}
\tag{3.1}
$$

(3.1)                            .

$$
\left.
\begin{array}{l}
0 \le a \le x_1 \\
0 \le a \le y_1 \\
(x_1 + y_1) - t \le a \le x_1 + y_1
\end{array}
\right\}
\tag{3.2}
$$

< 1> $t$, $x_1$, $y_1$                              ,

.                    $t$, $x_1$, $y_1$                         X   Y

$a$                        < 1>                    . < 1>

$\chi^2(1)$                              ,

. Cochran(1954)

Fisher              (Fisher's exact test)

,

.

$t$, $x_1$, $y_1$              ,        X   Y                $a$

,        $a$                          ,

,

.

## 3.2

< 1>                                   .

$$\chi^2 = \sum_{i,j} \frac{(E_{ij} - O_{ij})^2}{E_{ij}}$$

$$= \left( \frac{t}{x_1 y_1 (t - x_1)(t - y_1)} \right) t^2 a^2 - 2t x_1 y_1 a + (x_1 y_1)^2)$$

(3.3)

$a$        2              $\alpha_2$, 1              $\alpha_1$,          $\alpha_0$          (3.3)

.

$$\chi^2 = \alpha_2 a^2 + \alpha_1 a + \alpha_0$$

(3.4)

$$\alpha_2 = \frac{t^3}{x_1 y_1 (t - x_1)(t - y_1)}$$

$$\alpha_1 = \frac{- 2 t^2}{(t - x_1)(t - y_1)}$$

$$\alpha_0 = \frac{t x_1 y_1}{(t - x_1)(t - y_1)}$$

.  (3.4)        (3.2)                              .

$$\left.\begin{array}{l} \alpha_2 \geq 0 \\ \alpha_1 \leq 0 \\ \alpha_0 \geq 0 \end{array}\right\} \qquad (3.5)$$

< 1> .

$$H_0 : p_{ij} = p_{i \cdot} \, p_{\cdot j} \, (i = 1, 2 \, ; j = 1, 2)$$

$$H_1 : H_0 \qquad .$$

(3.6)

$$p_{ij} = \qquad x_i \qquad y_j$$

$$p_{i \cdot} = \qquad x_i$$

$$p_{\cdot j} = \qquad y_j$$

. $H_0$ $a$ .

$$\alpha_2 \, a^2 + \alpha_1 \, a + \alpha_0 \geq \chi_\alpha^2(1) \qquad (3.7)$$

$\chi_\alpha^2(1)$ $\alpha$ , 1 . $a_L$ $a$ ,
$a_U$ $a$ , $a_L$ $a_U$ .

$$a_L = \frac{-\alpha_1 - \sqrt{\alpha_1^2 - 4\alpha_2(\alpha_0 - \chi_\alpha^2(1))}}{2\alpha_2}$$

$$a_U = \frac{-\alpha_1 + \sqrt{\alpha_1^2 - 4\alpha_2(\alpha_0 - \chi_\alpha^2(1))}}{2\alpha_2}$$

(3.7) $a$ .

$$a \leq a_L , a \geq a_U \qquad (3.8)$$

X Y $H_0$ $a$
(3.5) 2 $\alpha_2 \geq 0$ , (3.8) $a$ (3.2)
.

## 3.3

< 2> < 1>
, , . $2 \times 2$ 8 , ,
, , $R : X \Rightarrow Y$ ,

,                     $Sup(X \Rightarrow Y)$, $Conf(X \Rightarrow Y)$, $Lift(X \Rightarrow Y)$                .

                              X    Y                                    $a$

          ,                              1                              ,                              $a$

                              1                                      .

                                                  $R : X \Rightarrow Y$        $Sup(X \Rightarrow Y) \geq min\_sup$

$Conf(X \Rightarrow Y) \geq min\_conf$                                            $a$  (

     )              $a$  (              )                    .

     $a$                      1              $a$

                    ,                                                  $a$                      1

          $a$

,                    X    Y                                                              (3.8)

     (3.2)              $a$

                              .

<    2> 2× 2                              ,          ,

| $X \Rightarrow Y$ | | $Y \Rightarrow X$ | |
|---|---|---|---|
| $Sup(X \Rightarrow Y) = \dfrac{a}{t}$ <br> $Conf(X \Rightarrow Y) = \dfrac{a}{x_1}$ <br> $Lift(X \Rightarrow Y) = \dfrac{ta}{x_1 y_1}$ | | $Sup(Y \Rightarrow X) = \dfrac{a}{t}$ <br> $Conf(Y \Rightarrow X) = \dfrac{a}{y_1}$ <br> $Lift(Y \Rightarrow X) = \dfrac{ta}{x_1 y_1}$ | |
| $X \Rightarrow \sim Y$ | | $\sim Y \Rightarrow X$ | |
| $Sup(X \Rightarrow \sim Y) = \dfrac{x_1 - a}{t}$ <br> $Conf(X \Rightarrow \sim Y) = \dfrac{x_1 - a}{x_1}$ <br> $Lift(X \Rightarrow \sim Y) = \dfrac{t(x_1 - a)}{x_1(t - y_1)}$ | | $Sup(\sim Y \Rightarrow X) = \dfrac{x_1 - a}{t}$ <br> $Conf(\sim Y \Rightarrow X) = \dfrac{x_1 - a}{t - y_1}$ <br> $Lift(\sim Y \Rightarrow X) = \dfrac{t(x_1 - a)}{x_1(t - y_1)}$ | |
| $\sim X \Rightarrow Y$ | | $Y \Rightarrow \sim X$ | |
| $Sup(\sim X \Rightarrow Y) = \dfrac{y_1 - a}{t}$ <br> $Conf(\sim X \Rightarrow Y) = \dfrac{y_1 - a}{t - x_1}$ <br> $Lift(\sim X \Rightarrow Y) = \dfrac{t(y_1 - a)}{y_1(t - x_1)}$ | | $Sup(Y \Rightarrow \sim X) = \dfrac{y_1 - a}{t}$ <br> $Conf(Y \Rightarrow \sim X) = \dfrac{y_1 - a}{y_1}$ <br> $Lift(Y \Rightarrow \sim X) = \dfrac{t(y_1 - a)}{y_1(t - x_1)}$ | |
| $\sim X \Rightarrow \sim Y$ | | $\sim Y \Rightarrow \sim X$ | |
| $Sup(\sim X \Rightarrow \sim Y) = \dfrac{t - (x_1 + y_1) + a}{t}$ <br> $Conf(\sim X \Rightarrow \sim Y) = \dfrac{t - (x_1 + y_{1)} + a}{t - x_1}$ <br> $Lift(\sim X \Rightarrow \sim Y) = \dfrac{t(t - (x_1 + y_1) + a)}{(t - x_1)(t - y_1)}$ | | $Sup(\sim Y \Rightarrow \sim X) = \dfrac{t - (x_1 + y_1) + a}{t}$ <br> $Conf(\sim Y \Rightarrow \sim X) = \dfrac{t - (x_1 + y_{1)} + a}{t - y_1}$ <br> $Lift(\sim Y \Rightarrow \sim X) = \dfrac{t(t - (x_1 + y_1) + a)}{(t - x_1)(t - y_1)}$ | |

* : not              .

3.4

(3.8)        (3.2)                          $a$                                                      <    2>

$Conf(X \Rightarrow Y)$                          .                    $Conf(X \Rightarrow Y)$        $Conf_{xy}$                          .

$$Conf_{xy} = \frac{a}{x_1} \quad , \qquad a \qquad\qquad X \quad Y$$

(3.8)        (3.2)                                          .

$a = x_1 Conf_{xy}$        (3.3)                          $Conf_{xy}$                                      .

$$\chi^2 = \frac{t}{x_1 y_1 (t - x_1)(t - y_1)} (t^2 x_1^2 \ Conf_{xy}^2 - 2tx_1 y_1 x_1 Conf_{xy} + x_1^2 y_1^2)$$

$$= \frac{tx_1}{y_1 (t - x_1)(t - y_1)} (t^2 \ Conf_{xy}^2 - 2ty_1 Conf_{xy} + y_1^2)$$

(3.9)

$Conf_{xy}$                2                          $\beta_2$, 1                          $\beta_1$,                  $\beta_0$                          (3.9)

.

$$\chi^2 = \beta_2 \ Conf_{xy}^2 + \beta_1 Conf_{xy} + \beta_0$$

(3.10)

$$\beta_2 = \frac{t^3 x_1}{y_1 (t - x_1)(t - y_1)}$$

$$\beta_1 = \frac{-2 \ t^2 x_1}{(t - x_1)(t - y_1)}$$

$$\beta_0 = \frac{ty_1}{(t - x_1)(t - y_1)}$$

.   (3.10)          (3.2)                                                              .

$$\beta_2 \geq 0, \ \beta_1 \leq 0, \ \beta_0 \geq 0$$

(3.11)

(3.10)      $X_\alpha^2(1)$

.

$$\beta_2 \ Conf_{xy}^2 + \beta_1 Conf_{xy} + \beta_0 \geq X_\alpha^2(1)$$

(3.12)

$\chi_\alpha^2(1)$                        $\alpha$        ,                    1                                      .

$Conf_{xy_L}$          $Conf_{xy}$                      ,        $Conf_{xy_U}$          $Conf_{xy}$                                          ,        $Conf_{xy_L}$

$Conf_{xy\,U}$ .

$$Conf_{xy\,L} = \frac{-\,\beta_1 - \sqrt{\beta_1^2 - 4\beta_2(\beta_0 - \chi_\alpha^2(1))}}{2\beta_2}$$

$$Conf_{xy\,U} = \frac{-\,\beta_1 + \sqrt{\beta_1^2 - 4\beta_2(\beta_0 - \chi_\alpha^2(1))}}{2\beta_2}$$

(3.12) $Conf_{xy}$ .

$$Conf_{xy} \le Conf_{xy\,L} \,, Conf_{xy} \ge Conf_{xy\,U} \qquad\qquad (3.13)$$

X    Y                                                               $Conf_{xy}$

(3.11)      2                   $\beta_2 \ge 0$     ,      (3.13)      $a$                              (3.2)

$Conf_{xy} = a/x_1$                              .

,

.

### 4.

3                                    ,

.                    $t,\ y_1,\ x_1$                          ,                                  $a$

,                                    ,                $t, y_1$                  $x_1,\ a$

,                                                                      .

X, Y                                  .

( $t$ )    100            ,            X

300        (1)              90            300        (0)                  10

.              Y                                          (1)            25

(0)              75              .            X    Y

,    300                                                $a$

.          <    3>          .

<    3>

|   |   | Y |   |   |
|---|---|---|---|---|
|   |   | 1 | 0 |   |
| X | 1 | $a$ | 90 - $a$ | 90 |
|   | 0 | 25 - $a$ | $a$ - 15 | 10 |
|   |   | 25 | 75 | 100 |

(3.2)    $t = 100, x_1 = 90, y_1 = 25$              $a$                                    ,

$$15 \leq a \leq 19 \tag{4.1}$$

(3.4)    $\alpha_2, \alpha_1, \alpha_0$                                    .

$$\alpha_2 = \frac{t^3}{x_1 y_1 (t - x_1)(t - y_1)} = \frac{100^3}{90 \times 25 \times (100 - 90) \times (100 - 25)} = 0.5926$$

$$\alpha_1 = \frac{-2\,t^2}{(t - x_1)(t - y_1)} = \frac{-2 \times 100^2}{(100 - 90)(100 - 25)} = -26.6667$$

$$\alpha_0 = \frac{t x_1 y_1}{(t - x_1)(t - y_1)} = \frac{100 \times 90 \times 25}{(100 - 90) \times (100 - 25)} = 300$$

$\alpha = 0.05$        $\chi^2(1) = 3.84146$    ,   $a$                                    .

$$a_L = \frac{-\alpha_1 - \sqrt{\alpha_1^2 - 4\alpha_2(\alpha_0 - \chi_\alpha^2(1))}}{2\alpha_2} = \frac{-(-26.6667) - \sqrt{19.095}}{2 \times 0.5926} = 19.955$$

$$a_U = \frac{-\alpha_1 + \sqrt{\alpha_1^2 - 4\alpha_2(\alpha_0 - \chi_\alpha^2(1))}}{2\alpha_2} = \frac{-(-26.5557) + \sqrt{9.095}}{2 \times 0.5926} = 25.044$$

$H_0$                    $a$

$$a \leq 19.955, a \geq 25.044 \tag{4.2}$$

(4.1)        (4.2)        $a$                                    .

$$15 \leq a \leq 19 \tag{4.3}$$

,              X    Y    $a$    $15 \leq a \leq 19$                                    .

(4.3)    <    1>                                    .

$$\frac{15}{100} \leq \mathrm{Sup}\,(X \Rightarrow Y) = \frac{a}{t} \leq \frac{19}{100}$$

$$\frac{15}{90} \leq \mathrm{Conf}(X \Rightarrow Y) = \frac{a}{x_1} \leq \frac{19}{90}$$

$$\frac{100 \times 15}{90 \times 25} \leq \mathrm{Lift}\,(X \Rightarrow Y) = \frac{ta}{x_1 y_1} \leq \frac{100 \times 19}{90 \times 25}$$

$R : X \Rightarrow Y$                                    $a = 15$          ,

$min\_sup = 15\%$        ,                            $min\_conf = 16.7\%$                    .

$Lift = 0.667$                    .              $<$    2$>$
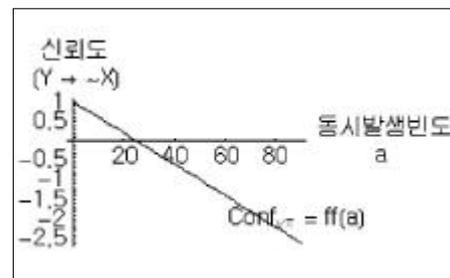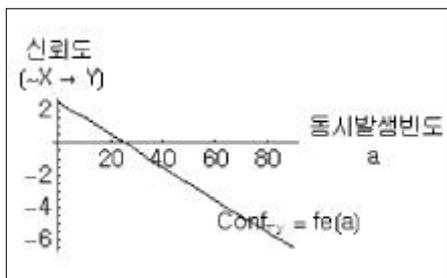
X       Y

$a = 15$         ,                $min\_conf = 16.7\%$                          .



$<$       1$>$



$<$      2$>$

$<$      1$>$   $<$     2$>$   $<$    3$>$                            $a$

.      $<$      3$>$   $<$    3$>$

$a$                                                        .

신뢰도
(X → Y)
$Conf_{xy} = fa(a)$
1
0.8
0.6
0.4
0.2
동시발생빈도
20 40 60 80
a

신뢰도
(Y → X)
$Conf_{yx} = fb(a)$
3.5
3
2.5
2
1.5
1
0.5
동시발생빈도
20 40 60 80
a

신뢰도
(X → ~Y)
1
0.8
0.6
0.4
0.2
$Conf_{x\overline{y}} = fc(a)$
동시발생빈도
20 40 60 80
a

신뢰도
(~Y → X)
1.2
1
0.8
0.6
0.4
0.2
$Conf_{\overline{y}x} = fd(a)$
동시발생빈도
20 40 60 80
a

신뢰도
(~X → Y)
2
동시발생빈도
20 40 60 80
a
-2
-4
-6
$Conf_{\overline{x}y} = fe(a)$

신뢰도
(Y → ~X)
1
0.5
동시발생빈도
20 40 60 80
a
-0.5
-1
-1.5
-2
-2.5
$Conf_{y\overline{x}} = ff(a)$

신뢰도
(~X → ~Y)
$Conf_{\overline{x}\overline{y}} = fg(a)$
6
4
2
동시발생빈도
20 40 60 80
a

신뢰도
(~Y → ~X)
$Conf_{\overline{y}\overline{x}} = fh(a)$
1
0.8
0.6
0.4
0.2
동시발생빈도
20 40 60 80
a
-0.2

< 3>

< 4> < 3>                                                   *a*

,                                              (%)                    .

&lt;    4&gt;                                                                          (        )

5.


,                          ,

,


.



.



,


.

[1]          (2000).                                                            , *J.Basic Sci. Res. Kyungsan Univ*. 4(1), P 91-98.

[2] Han, J., Kamber M. (2001). Data Mining : *Concepts and Techniques*, Morgan Kaufmann Publishers.

[3] Reynolds, H. T. (1977). *Analysis of Nominal Data, Quantitative Applications in the Social Sciences,* Series No. 7, Sage Publications.

[4] Stephen, E. F. (1976). *The Analysis of Cross-Classfied Categorical Data*, Second Edition.

[5] Cochran, W. G. (1954), Some methods for strengthening the common $\chi^2$ tests, *Biometrics*, 10.

[6] Agrawal, R., Imielinski, R., Swami, A. (1993). Mining association rules between sets of items in

large databases, *Proceedings of the ACM SIGMOD Conference on Management of Data*, Washington, D.C.

[7] Agrawal, R., Srikant, R. (1994). Fast algorithms for mining association rules, *Proceedings of the 20th VLDB Conference*, Santiago, Chile.

[8] Bing, L., Wynne, H., Yiming, M. (1999). Mining Association Rules with Multiple minimum Supports, *Proceedings of ACM KDD*-99.

[9] Toivonen, H. (1996). Sampling Large Database for Association Rules, *Proceedings of the 22nd VLDB Conference Mumbai(Bombay)*, India.

[10] Cheung, D. W., Han, J., Ng, V., Wong, C. Y. (1996). Maintenance of discovered association rules in large databases: An incremental updating technique, *Int'l Conference on Data Engineering*, New Orleans, Louisiana.

[11] Sergey, B., Rajeev, M., Jeffrey, D. U., Shalom, T. (1997). Dynamic itemset counting and implication rules for market data, *Proceedings of ACM SIGMOD Conference on Management of Data*.

[12] Silverstein, C., Brin, S., Motwani, R. (1998). Beyond Market Baskets: Generalizing Association Rules to Dependence Rules, *Data Mining and Knowledge Discovery*, No.2, P 39-68.

[13] Park, J. S., Chen, M. S., and Philip, S.Y . (1995). An effective hash-based algorithms for mining association rules, *Proceedings of ACM SIGMOD Conference on Management of Data*.

[14] Agrawal, R., John, C. S. (1996). Parallel Mining of Association Rules, *IEEE Transactions on Knowledge and Data Engineering*, Vol.8, No.6.

[15] Cheung, D. W., Han, J., Ng, V., Fu, A. W., Fu, Y. (1996). A Fast distribution algorithm for mining association rules, *Int's Conference on Parallel and Distributes Information System*, Miami Beach, Florida.

[16] Han, J., Pei, J. (2000). Mining Frequent Patterns by Pattern-Growth : Methodology and Implications, *SIGKDD Explorations*.

[17] Hipp, J., Guntzer, U., Gholamreza (2000). Algorithms for Association Rule Mining - A General Survey and Comparison, *SIGKDD Explorations,* Vol. 2, Issue 1, P 58 - 64.

[18] Saygin, Y., Vassilios, S. V., Clifton, C. (2002). Using Unknowns to Prevent Discovery of Association Rules, *2002 Conference on Research Issues in Data Engineering*.

[19] John, F. R., Rice, S. (2001). What's Interesting About Cricket? - On Thresholds and Anticipation in Discovered Rules, *SIGKDD Explorations,* Vol.3 Issue 1, P 1-5.

[20] Khalil, M. A., Nagwa, M. E, Taha, Y. (2000). *A note on 'Beyond Market Baskets: Generalizing Association Rules to Correlations'*.