

Statistical Decision making of Association Threshold in Association Rule Data Mining

Hee Chang Park¹⁾ · Geum Min Song²⁾

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 calculate the range of the value that two item sets are occurred simultaneously, and find the minimum confidence threshold values.

Keywords : , , ,

1.

(data mining) (mine)
가 가 . ,
가 가 .
(association rule)
, (decision tree), (neural network), (clustering),
(genetic algorithm), (bayesian network), -
(memory-based reasoning)
(support), (confidence), (lift)
, 가

-
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.

Agrawal (1993) , Agrawal (1994)
 가
 Apriori AprioriTid
 . Park (1995)
 , Toivonen(1996) partitioning
 가
 가 sampling
 . Cheung (1996)
 FUP(fast update) , Sergey (1997) 가
 itemset counting) . Liu (1999) DIC(dynamic
 가 DHP(direct hashing and pruning) . Saygin (2002)
 .
 가 Silverstein (1997) Fast
 (0)
 (1)
 (minimum confidence) (minimum support),
 가
 ,
 가
 ,
 2 가
 , 3 ,
 , 4 가
 . 5 가
 .
2.

가

- (transaction) : 가
- (item set) : ,
- (candidate item set) :
- (frequent item set) : 가

k- 가 $\mathcal{I} = \{i_1, i_2, \dots, i_k\}$ $T \subseteq \mathcal{I}$,
 T TID 가 $A \subseteq T$
 T A

$R: X \Rightarrow Y$

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

: $Conf(X \Rightarrow Y) = P(Y|X) = \frac{P(X \cap Y)}{P(X)}$

: $Lift(X \Rightarrow Y) = \frac{P(Y|X)}{P(Y)} = \frac{P(X \cap Y)}{P(X)P(Y)}$

$X \subset T, Y \subset T, \text{ and } X \cap Y \neq \emptyset$

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

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

가 (Lift)
가

Apriori k- k-1
(join step) (candidate k-itemsets) 가 가
가 가 (prune
step) 가

3.

3.1 가

(association)

(test of independence)

2 × 2 (contingency table)

		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
 $x_0 = t - x_1 :$ X가
 $y_0 = t - y_1 :$ Y가
 (3.1)

$$\left. \begin{aligned}
 0 \leq a \leq t \\
 0 \leq x_1 - a \leq t \\
 0 \leq y_1 - a \leq t \\
 0 \leq t - (x_1 + y_1) + a \leq t \\
 0 \leq a \leq x_1 \\
 0 \leq a \leq y_1
 \end{aligned} \right\} \quad (3.1)$$

(3.1)

$$\max(0, (x_1 + y_1) - t) \leq a \leq \min(x_1, y_1) \quad (3.2)$$

$\langle 1 \rangle$ t, x_1, y_1 가 ,
 t, x_1, y_1
 X Y a $\langle 1 \rangle$
 $\langle 1 \rangle$ 가 $\chi^2(1)$

Fisher . Cochran(1954)
(Fisher's exact test) , 가 가

t, x_1, y_1 , X Y a
, a

, 가

3.2

< 1>

$$\chi^2 = \sum_{i,j} \frac{(E_{ij} - O_{ij})^2}{E_{ij}} \tag{3.3}$$

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

a^2 $\alpha_2, 1$ $\alpha_1,$ α_0 (3.3)

$$\chi^2 = \alpha_2 a^2 + \alpha_1 a + \alpha_0 \tag{3.4}$$

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

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

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

(3.4) (3.2)

$$\left. \begin{matrix} \alpha_2 \geq 0 \\ \alpha_1 \leq 0 \\ \alpha_0 \geq 0 \end{matrix} \right\} \tag{3.5}$$

< 1>

가 가

$$\begin{aligned}
 H_0 &: p_{ij} = p_{i \cdot} p_{\cdot j} \quad (i = 1, 2; j = 1, 2) \\
 H_1 &: H_0 \text{가 아니다}
 \end{aligned}
 \tag{3.6}$$

$$\begin{aligned}
 p_{ij} &= \frac{x_i y_j}{a} \\
 p_{i \cdot} &= \frac{x_i}{a} \\
 p_{\cdot j} &= \frac{y_j}{a}
 \end{aligned}$$

가 H_0 를 만족하는 a 가 존재하는지 판별하기 위하여

$$\alpha_2 a^2 + \alpha_1 a + \alpha_0 \geq \chi_\alpha^2(1)
 \tag{3.7}$$

$\chi_\alpha^2(1)$ 를 α 로 나타내면, $\alpha > 1$ 이므로 a_L 와 a_U 가 존재한다.

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

(3.7)을 만족하는 a 는 $a_L \leq a \leq a_U$ 이다.

$$a \leq a_L, a \geq a_U
 \tag{3.8}$$

a 가 존재하는지 판별하기 위하여 (3.5)에서 $\alpha_2 \geq 0$ 인 경우, (3.8)에서 a 가 존재하는지 판별한다.

3.3

$\langle 2 \rangle < 1 \rangle$ 인 경우, 2×2 contingency table의 경우, $R: X \Rightarrow Y$ 의 $\text{Sup}(X \Rightarrow Y)$, $\text{Conf}(X \Rightarrow Y)$, $\text{Lift}(X \Rightarrow Y)$ 를 계산한다.

X 와 Y 가 독립이면, a 가 존재한다.

$$\begin{aligned}
 & R : X \Rightarrow Y \quad \text{Sup}(X \Rightarrow Y) \geq \text{min_sup} \\
 & \text{Conf}(X \Rightarrow Y) \geq \text{min_conf} \\
 & a (\quad) \quad a (\quad) \quad \text{가 } a \quad \text{가 } 1 \quad a \\
 & \quad \quad \quad \text{가 } a \quad \quad \quad 1 \quad a \\
 & \quad \quad \quad \text{가 } \quad \quad \quad \text{가 } \quad \quad \quad \text{가 } \quad \quad \quad \text{가 } \quad \quad \quad (3.8) \\
 & X \quad Y \text{가} \\
 & (3.2) \quad a \quad \text{가} \quad \quad \quad \text{가} \quad \quad \quad \text{가}
 \end{aligned}$$

< 2 > 2 × 2

$X \Rightarrow Y$	$Y \Rightarrow X$
$\text{Sup}(X \Rightarrow Y) = \frac{a}{t}$ $\text{Conf}(X \Rightarrow Y) = \frac{a}{x_1}$ $\text{L if } t(X \Rightarrow Y) = \frac{ta}{x_1 y_1}$	$\text{Sup}(Y \Rightarrow X) = \frac{a}{t}$ $\text{Conf}(Y \Rightarrow X) = \frac{a}{y_1}$ $\text{L if } t(Y \Rightarrow X) = \frac{ta}{x_1 y_1}$
$X \Rightarrow \sim Y$	$\sim Y \Rightarrow X$
$\text{Sup}(X \Rightarrow \sim Y) = \frac{x_1 - a}{t}$ $\text{Conf}(X \Rightarrow \sim Y) = \frac{x_1 - a}{x_1}$ $\text{L if } t(X \Rightarrow \sim Y) = \frac{t(x_1 - a)}{x_1(t - y_1)}$	$\text{Sup}(\sim Y \Rightarrow X) = \frac{x_1 - a}{t}$ $\text{Conf}(\sim Y \Rightarrow X) = \frac{x_1 - a}{t - y_1}$ $\text{L if } t(\sim Y \Rightarrow X) = \frac{t(x_1 - a)}{x_1(t - y_1)}$
$\sim X \Rightarrow Y$	$Y \Rightarrow \sim X$
$\text{Sup}(\sim X \Rightarrow Y) = \frac{y_1 - a}{t}$ $\text{Conf}(\sim X \Rightarrow Y) = \frac{y_1 - a}{t - x_1}$ $\text{L if } t(\sim X \Rightarrow Y) = \frac{t(y_1 - a)}{y_1(t - x_1)}$	$\text{Sup}(Y \Rightarrow \sim X) = \frac{y_1 - a}{t}$ $\text{Conf}(Y \Rightarrow \sim X) = \frac{y_1 - a}{y_1}$ $\text{L if } t(Y \Rightarrow \sim X) = \frac{t(y_1 - a)}{y_1(t - x_1)}$
$\sim X \Rightarrow \sim Y$	$\sim Y \Rightarrow \sim X$
$\text{Sup}(\sim X \Rightarrow \sim Y) = \frac{t - (x_1 + y_1) + a}{t}$ $\text{Conf}(\sim X \Rightarrow \sim Y) = \frac{t - (x_1 + y_1) + a}{t - x_1}$ $\text{L if } t(\sim X \Rightarrow \sim Y) = \frac{t(t - (x_1 + y_1) + a)}{(t - x_1)(t - y_1)}$	$\text{Sup}(\sim Y \Rightarrow \sim X) = \frac{t - (x_1 + y_1) + a}{t}$ $\text{Conf}(\sim Y \Rightarrow \sim X) = \frac{t - (x_1 + y_1) + a}{t - y_1}$ $\text{L if } t(\sim Y \Rightarrow \sim X) = \frac{t(t - (x_1 + y_1) + a)}{(t - x_1)(t - y_1)}$

* \sim : not

3.4

(3.8) (3.2) a $\langle 2 \rangle$
 $Conf(X \Rightarrow Y)$ $Conf(X \Rightarrow Y)$
 $Conf_{xy}$ $Conf_{xy} = \frac{a}{x_1}$, a X
 Y가 (3.8) (3.2)
 $a = x_1 Conf_{xy}$ (3.3) $Conf_{xy}$

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

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

(3.9)

$Conf_{xy}$ 2 $\beta_2, 1$ β_1, β_0
 (3.9)

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

(3.10)

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

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

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

(3.10) (3.2)

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

(3.11)

가 (3.10) $X_\alpha^2(1)$

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

(3.12)

$\chi_\alpha^2(1)$ α , 가 1

$Conf_{xyL}$ $Conf_{xy}$, $Conf_{xyU}$ $Conf_{xy}$, $Conf_{xyL}$
 $Conf_{xyU}$.

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

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

(3.12) $Conf_{xy}$.

$Conf_{xy} \leq Conf_{xyL}$, $Conf_{xy} \geq Conf_{xyU}$ (3.13)

$Conf_{xy}$ X Y가 (3.11) 2 $\beta_2 \geq 0$ 가 (3.13) a
(3.2) $Conf_{xy} = a/x_1$ 가 .
, 가

4. ()

3 ,
가 t, y_1, x_1 ,
가 a ,
, t, y_1 x_1, a ,
가 .
X, Y 가 X
(t) 100 ,
(1) 90 X 300
(0) 300 10 Y
(1) 25 .
(0) 75 X Y가 ,
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 가

$$15 \leq a \leq 25 \tag{4.1}$$

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

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

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

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

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

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

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

가 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 가 $15 \leq a \leq 19$
 (4.3) $< 1 >$.

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

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

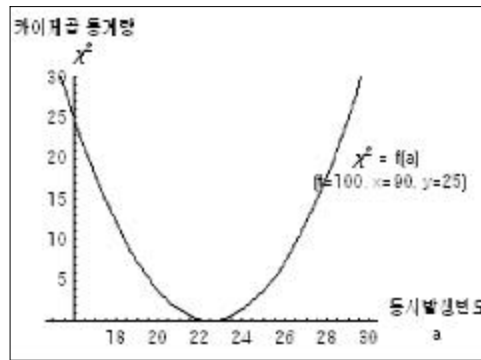
$$\frac{100 \times 15}{90 \times 25} \leq \text{Lift}(X \Rightarrow Y) = \frac{ta}{x_1 y_1} \leq \frac{100 \times 19}{90 \times 25}$$

$R: X \Rightarrow Y$ 가

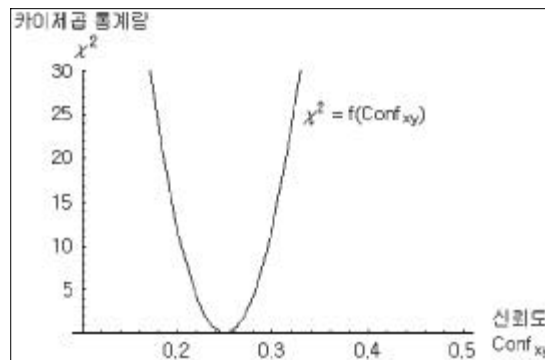
$a = 15$,
 $\text{min_conf} = 16.7\%$
 < 2>
 $X \Rightarrow Y$
 $\text{min_conf} = 16.7\%$

$\text{min_sup} = 15\%$,
 $\text{Lift} = 0.667$

가 $a = 15$,

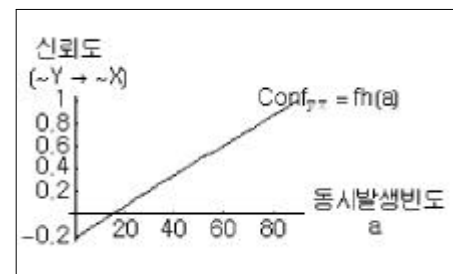
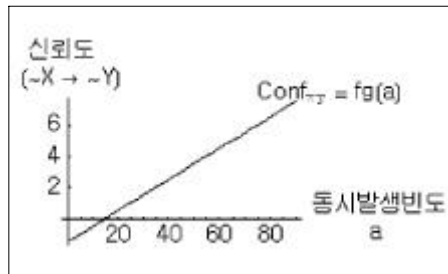
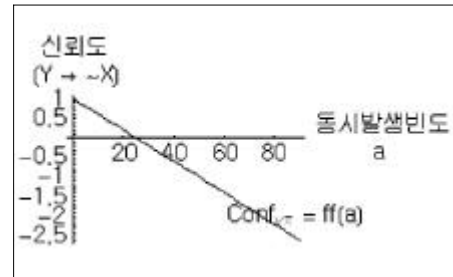
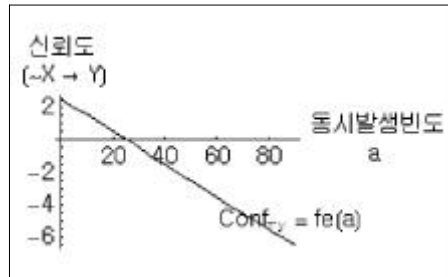
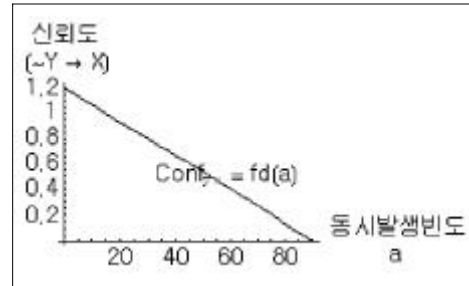
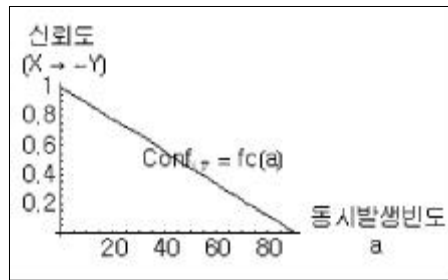
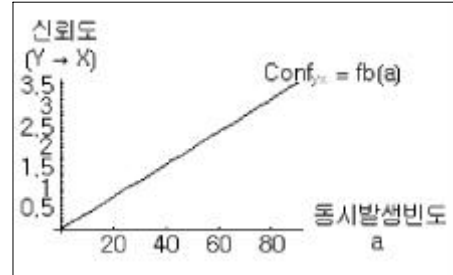
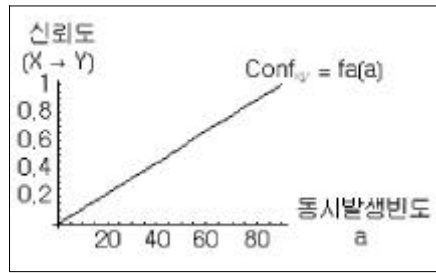


< 1>



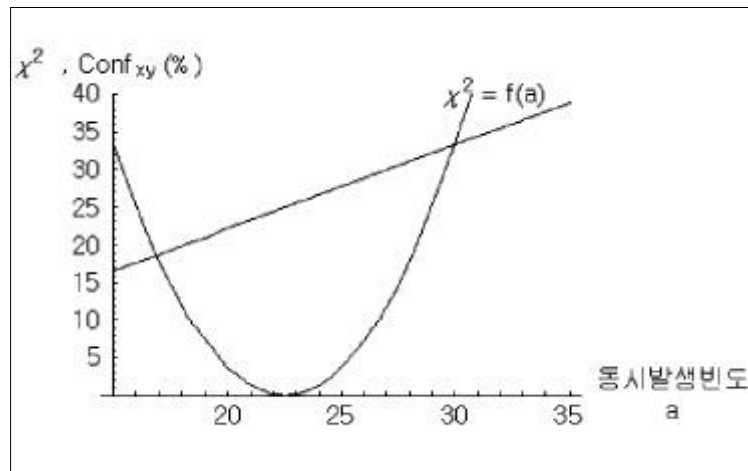
< 2>

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



< 3>

< 4> < 3> a (%)



< 4> ()

5.

가 , 가 가

가

가 가

1. 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.
2. Agrawal, R., John, C.S. (1996). Parallel Mining of Association Rules, *IEEE Transactions on Knowledge and Data Engineering*, Vol.8, No.6.
 3. Agrawal, R., Srikant, R. (1994). Fast algorithms for mining association rules, *Proceedings of the 20th VLDB Conference*, Santiago, Chile.
 4. Bing, L., Wynne, H., Yiming, M. (1999). Mining Association Rules with Multiple minimum Supports, *Proceedings of A CM KDD-99*.
 5. 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.
 6. Cheung, D.W., Han, J., Ng, V., Fu, A.W., Fu, Y. (1996). A Fast distribution algorithm for mining association rules, *Int'l Conference on Parallel and Distributed Information System*, Miami Beach, Florida.
 7. Cochran, W.G. (1954), *Some methods for strengthening the common χ^2 tests*, *Biometrics*, 10.
 8. Han, J., Kamber M. (2001). *Data Mining : Concepts and Techniques*, Morgan Kaufmann Publishers.
 9. Han, J., Pei, J. (2000). Mining Frequent Patterns by Pattern-Growth : Methodology and Implications, *SIGKDD Explorations*.
 10. Park, J.S., Chen, M.S., and Philip, S.Y. (1995). An effective hash-based algorithms for mining association rules, *Proceedings of A CM SIGMOD Conference on Management of Data*.
 11. Saygin, Y., Vassilios, S.V., Clifton, C. (2002). Using Unknowns to Prevent Discovery of Association Rules, *2002 Conference on Research Issues in Data Engineering*.
 12. Sergey, B., Rajeev, M., Jeffrey, D.U., Shalom, T. (1997). Dynamic itemset counting and implication rules for market data, *Proceedings of A CM SIGMOD Conference on Management of Data*.
 13. 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.
 14. Toivonen, H. (1996). Sampling Large Database for Association Rules, *Proceedings of the 22nd VLDB Conference*, Mumbai(Bombay), India.