

Statistical Independence from the Viewpoint of Linear Algebra

Shusaku Tsumoto

Department of Medical Informatics,
Shimane University, School of Medicine,
89-1 Enya-cho Izumo 693-8501 Japan
tsumoto@computer.org

Abstract. A contingency table summarizes the conditional frequencies of two attributes and shows how these two attributes are dependent on each other with the information on a partition of universe generated by these attributes. Thus, this table can be viewed as a relation between two attributes with respect to information granularity. This paper focuses on statistical independence in a contingency table from the viewpoint of granular computing, which shows that statistical independence in a contingency table is a special form of linear dependence. The discussions also show that when a contingency table is viewed as a matrix, its rank is equal to 1.0. Thus, the degree of independence, rank plays a very important role in extracting a probabilistic model from a given contingency table.

1 Introduction

Statistical independence between two attributes is a very important concept in data mining and statistics. The definition $P(A, B) = P(A)P(B)$ show that the joint probability of A and B is the product of both probabilities. This gives several useful formula, such as $P(A|B) = P(A)$, $P(B|A) = P(B)$. In a data mining context, these formulae show that these two attributes may not be correlated with each other. Thus, when A or B is a classification target, the other attribute may not play an important role in its classification.

Although independence is a very important concept, it has not been fully and formally investigated as a relation between two attributes.

In this paper, a statistical independence in a contingency table is focused on from the viewpoint of granular computing.

The first important observation is that a contingency table compares two attributes with respect to information granularity. It is shown from the definition that statistical independence in a contingency table is a special form of linear dependence of two attributes. Especially, when the table is viewed as a matrix, the above discussion shows that the rank of the matrix is equal to 1.0. Also, the results also show that partial statistical independence can be observed.

The second important observation is that matrix algebra is a key point of analysis of this table. A contingency table can be viewed as a matrix and several

operations and ideas of matrix theory are introduced into the analysis of the contingency table.

The paper is organized as follows: Section 2 discusses the characteristics of contingency tables. Section 3 shows the conditions on statistical independence for a 2×2 table. Section 4 gives those for a $2 \times n$ table. Section 5 extends these results into a multi-way contingency table. Section 6 discusses statistical independence from matrix theory. Finally, Section 7 concludes this paper.

2 Contingency Table from Rough Sets

2.1 Rough Sets Notations

In the subsequent sections, the following notations is adopted, which is introduced in [7]. Let U denote a nonempty, finite set called the universe and A denote a nonempty, finite set of attributes, i.e., $a : U \rightarrow V_a$ for $a \in A$, where V_a is called the domain of a , respectively. Then, a decision table is defined as an information system, $A = (U, A \cup \{\mathcal{D}\})$, where $\{\mathcal{D}\}$ is a set of given decision attributes. The atomic formulas over $B \subseteq A \cup \{\mathcal{D}\}$ and V are expressions of the form $[a = v]$, called descriptors over B , where $a \in B$ and $v \in V_a$. The set $F(B, V)$ of formulas over B is the least set containing all atomic formulas over B and closed with respect to disjunction, conjunction and negation. For each $f \in F(B, V)$, f_A denote the meaning of f in A , i.e., the set of all objects in U with property f , defined inductively as follows.

1. If f is of the form $[a = v]$ then, $f_A = \{s \in U | a(s) = v\}$
2. $(f \wedge g)_A = f_A \cap g_A$; $(f \vee g)_A = f_A \cup g_A$; $(\neg f)_A = U - f_A$

By using this framework, classification accuracy and coverage, or true positive rate is defined as follows.

Definition 1.

Let R and D denote a formula in $F(B, V)$ and a set of objects whose decision attribute is given as \lceil , respectively. Classification accuracy and coverage(true positive rate) for $R \rightarrow \mathcal{D}$ is defined as:

$$\alpha_R(D) = \frac{|R_A \cap D|}{|R_A|} (= P(D|R)), \text{ and } \kappa_R(D) = \frac{|R_A \cap D|}{|D|} (= P(R|D)),$$

where $|A|$ denotes the cardinality of a set A , $\alpha_R(D)$ denotes a classification accuracy of R as to classification of \mathcal{D} , and $\kappa_R(D)$ denotes a coverage, or a true positive rate of R to \mathcal{D} , respectively.

2.2 Two-Way Contingency Table

From the viewpoint of information systems, a contingency table summarizes the relation between two attributes with respect to frequencies. This viewpoint has

already been discussed in [10, 11]. However, this study focuses on more statistical interpretation of this table.

Definition 2. Let R_1 and R_2 denote binary attributes in an attribute space A . A contingency table is a table of a set of the meaning of the following formulas: $|[R_1 = 0]_A|, |[R_1 = 1]_A|, |[R_2 = 0]_A|, |[R_1 = 1]_A|, |[R_1 = 0 \wedge R_2 = 0]_A|, |[R_1 = 0 \wedge R_2 = 1]_A|, |[R_1 = 1 \wedge R_2 = 0]_A|, |[R_1 = 1 \wedge R_2 = 1]_A|, |[R_1 = 0 \vee R_1 = 1]_A| (= |U|)$. This table is arranged into the form shown in Table 1, where: $|[R_1 = 0]_A| = x_{11} + x_{21} = x_{.1}$, $|[R_1 = 1]_A| = x_{12} + x_{22} = x_{.2}$, $|[R_2 = 0]_A| = x_{11} + x_{12} = x_{1.}$, $|[R_2 = 1]_A| = x_{21} + x_{22} = x_{2.}$, $|[R_1 = 0 \wedge R_2 = 0]_A| = x_{11}$, $|[R_1 = 0 \wedge R_2 = 1]_A| = x_{21}$, $|[R_1 = 1 \wedge R_2 = 0]_A| = x_{12}$, $|[R_1 = 1 \wedge R_2 = 1]_A| = x_{22}$, $|[R_1 = 0 \vee R_1 = 1]_A| = x_{.1} + x_{.2} = x_{..} (= |U|)$.

Table 1. Two way Contingency Table

	$R_1 = 0$	$R_1 = 1$	
$R_2 = 0$	x_{11}	x_{12}	$x_{1.}$
$R_2 = 1$	x_{21}	x_{22}	$x_{2.}$
	$x_{.1}$	$x_{.2}$	$x_{..}$
			$(= U = N)$

From this table, accuracy and coverage for $[R_1 = 0] \rightarrow [R_2 = 0]$ are defined as:

$$\alpha_{[R_1=0]}([R_2 = 0]) = \frac{|[R_1 = 0 \wedge R_2 = 0]_A|}{|[R_1 = 0]_A|} = \frac{x_{11}}{x_{.1}},$$

and

$$\kappa_{[R_1=0]}([R_2 = 0]) = \frac{|[R_1 = 0 \wedge R_2 = 0]_A|}{|[R_2 = 0]_A|} = \frac{x_{11}}{x_{1.}}.$$

2.3 Multi-way Contingency Table

Two-way contingency table can be extended into a contingency table for multi-nominal attributes.

Definition 3. Let R_1 and R_2 denote multinominal attributes in an attribute space A which have m and n values. A contingency tables is a table of a set of the meaning of the following formulas: $|[R_1 = A_j]_A|, |[R_2 = B_i]_A|, |[R_1 = A_j \wedge R_2 = B_i]_A|, |[R_1 = A_1 \wedge R_1 = A_2 \wedge \dots \wedge R_1 = A_m]_A|, |[R_2 = B_1 \wedge R_2 = A_2 \wedge \dots \wedge R_2 = A_n]_A|$ and $|U|$ ($i = 1, 2, 3, \dots, n$ and $j = 1, 2, 3, \dots, m$). This table is arranged into the form shown in Table 1, where: $|[R_1 = A_j]_A| = \sum_{i=1}^m x_{1i} = x_{.j}$, $|[R_2 = B_i]_A| = \sum_{j=1}^n x_{ji} = x_{i.}$, $|[R_1 = A_j \wedge R_2 = B_i]_A| = x_{ij}$, $|U| = N = x_{..}$ ($i = 1, 2, 3, \dots, n$ and $j = 1, 2, 3, \dots, m$).

Table 2. Contingency Table ($m \times n$)

	A_1	A_2	\cdots	A_n	Sum
B_1	x_{11}	x_{12}	\cdots	x_{1n}	$x_{1\cdot}$
B_2	x_{21}	x_{22}	\cdots	x_{2n}	$x_{2\cdot}$
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
B_m	x_{m1}	x_{m2}	\cdots	x_{mn}	$x_{m\cdot}$
Sum	$x_{\cdot 1}$	$x_{\cdot 2}$	\cdots	$x_{\cdot n}$	$x_{\cdot\cdot} = U = N$

3 Statistical Independence in 2×2 Contingency Table

Let us consider a contingency table shown in Table 1. Statistical independence between R_1 and R_2 gives:

$$\begin{aligned}
 P([R_1 = 0], [R_2 = 0]) &= P([R_1 = 0]) \times P([R_2 = 0]) \\
 P([R_1 = 0], [R_2 = 1]) &= P([R_1 = 0]) \times P([R_2 = 1]) \\
 P([R_1 = 1], [R_2 = 0]) &= P([R_1 = 1]) \times P([R_2 = 0]) \\
 P([R_1 = 1], [R_2 = 1]) &= P([R_1 = 1]) \times P([R_2 = 1])
 \end{aligned}$$

Since each probability is given as a ratio of each cell to N , the above equations are calculated as:

$$\begin{aligned}
 \frac{x_{11}}{N} &= \frac{x_{11} + x_{12}}{N} \times \frac{x_{11} + x_{21}}{N} \\
 \frac{x_{12}}{N} &= \frac{x_{11} + x_{12}}{N} \times \frac{x_{12} + x_{22}}{N} \\
 \frac{x_{21}}{N} &= \frac{x_{21} + x_{22}}{N} \times \frac{x_{11} + x_{21}}{N} \\
 \frac{x_{22}}{N} &= \frac{x_{21} + x_{22}}{N} \times \frac{x_{12} + x_{22}}{N}
 \end{aligned}$$

Since $N = \sum_{i,j} x_{ij}$, the following formula will be obtained from these four formulae.

$$x_{11}x_{22} = x_{12}x_{21} \text{ or } x_{11}x_{22} - x_{12}x_{21} = 0$$

Thus,

Theorem 1. *If two attributes in a contingency table shown in Table 1 are statistical independent, the following equation holds:*

$$x_{11}x_{22} - x_{12}x_{21} = 0 \tag{1}$$

□

It is notable that the above equation corresponds to the fact that the determinant of a matrix corresponding to this table is equal to 0. Also, when these four values are not equal to 0, the equation 1 can be transformed into:

$$\frac{x_{11}}{x_{21}} = \frac{x_{12}}{x_{22}}.$$

Let us assume that the above ratio is equal to $C(\text{constant})$. Then, since $x_{11} = Cx_{21}$ and $x_{12} = Cx_{22}$, the following equation is obtained.

$$\frac{x_{11} + x_{12}}{x_{21} + x_{22}} = \frac{C(x_{21} + x_{22})}{x_{21} + x_{22}} = C = \frac{x_{11}}{x_{21}} = \frac{x_{12}}{x_{22}}. \tag{2}$$

It is notable that this discussion can be easily extended into a $2 \times n$ contingency table where $n > 3$. The important equation will be extended into

$$\frac{x_{11}}{x_{21}} = \frac{x_{12}}{x_{22}} = \dots = \frac{x_{1n}}{x_{2n}} = \frac{x_{11} + x_{12} + \dots + x_{1n}}{x_{21} + x_{22} + \dots + x_{2n}} = \frac{\sum_{k=1}^n x_{1k}}{\sum_{k=1}^n x_{2k}} \tag{3}$$

Thus,

Theorem 2. *If two attributes in a contingency table ($2 \times k(k = 2, \dots, n)$) are statistical independent, the following equations hold:*

$$x_{11}x_{22} - x_{12}x_{21} = x_{12}x_{23} - x_{13}x_{22} = \dots = x_{1n}x_{21} - x_{11}x_{2n} = 0 \tag{4}$$

□

It is also notable that this equation is the same as the equation on collinearity of projective geometry [2].

4 Statistical Independence in $m \times n$ Contingency Table

Let us consider a $m \times n$ contingency table shown in Table 2. Statistical independence of R_1 and R_2 gives the following formulae:

$$P([R_1 = A_i, R_2 = B_j]) = P([R_1 = A_i])P([R_2 = B_j])$$

$(i = 1, \dots, m, j = 1, \dots, n).$

According to the definition of the table,

$$\frac{x_{ij}}{N} = \frac{\sum_{k=1}^n x_{ik}}{N} \times \frac{\sum_{l=1}^m x_{lj}}{N}. \tag{5}$$

Thus, we have obtained:

$$x_{ij} = \frac{\sum_{k=1}^n x_{ik} \times \sum_{l=1}^m x_{lj}}{N}. \tag{6}$$

Thus, for a fixed j ,

$$\frac{x_{i_a j}}{x_{i_b j}} = \frac{\sum_{k=1}^n x_{i_a k}}{\sum_{k=1}^n x_{i_b k}}$$

In the same way, for a fixed i ,

$$\frac{x_{ij_a}}{x_{ij_b}} = \frac{\sum_{l=1}^m x_{lj_a}}{\sum_{l=1}^m x_{lj_b}}$$

Since this relation will hold for any j , the following equation is obtained:

$$\frac{x_{i_a 1}}{x_{i_b 1}} = \frac{x_{i_a 2}}{x_{i_b 2}} \dots = \frac{x_{i_a n}}{x_{i_b n}} = \frac{\sum_{k=1}^n x_{i_a k}}{\sum_{k=1}^n x_{i_b k}}. \quad (7)$$

Since the right hand side of the above equation will be constant, thus all the ratios are constant. Thus,

Theorem 3. *If two attributes in a contingency table shown in Table 2 are statistical independent, the following equations hold:*

$$\frac{x_{i_a 1}}{x_{i_b 1}} = \frac{x_{i_a 2}}{x_{i_b 2}} \dots = \frac{x_{i_a n}}{x_{i_b n}} = \text{const.} \quad (8)$$

for all rows: i_a and i_b ($i_a, i_b = 1, 2, \dots, m$).

□

5 Contingency Matrix

The meaning of the above discussions will become much clearer when we view a contingency table as a matrix.

Definition 4. *A corresponding matrix $C_{T_{a,b}}$ is defined as a matrix the element of which are equal to the value of the corresponding contingency table $T_{a,b}$ of two attributes a and b , except for marginal values.*

Definition 5. *The rank of a table is defined as the rank of its corresponding matrix. The maximum value of the rank is equal to the size of (square) matrix, denoted by r .*

The contingency matrix of Table 2($T(R_1, R_2)$) is defined as $C_{T_{R_1, R_2}}$ as below:

$$\begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix}$$

5.1 Independence of 2×2 Contingency Table

The results in Section 3 corresponds to the degree of independence in matrix theory. Let us assume that a contingency table is given as Table 1. Then the corresponding matrix ($C_{T_{R_1, R_2}}$) is given as:

$$\begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix},$$

Then,

Proposition 1. *The determinant of $\det(C_{T_{R_1, R_2}})$ is equal to $x_{11}x_{22} - x_{12}x_{21}$,*

Proposition 2. *The rank will be:*

$$\text{rank} = \begin{cases} 2, & \text{if } \det(C_{T_{R_1, R_2}}) \neq 0 \\ 1, & \text{if } \det(C_{T_{R_1, R_2}}) = 0 \end{cases}$$

From Theorem 1,

Theorem 4. *If the rank of the corresponding matrix of a 2times2 contingency table is 1, then two attributes in a given contingency table are statistically independent. Thus,*

$$\text{rank} = \begin{cases} 2, & \text{dependent} \\ 1, & \text{statistical independent} \end{cases}$$

This discussion can be extended into $2 \times n$ tables.

Theorem 5. *If the rank of the corresponding matrix of a $2 \times n$ contingency table is 1, then two attributes in a given contingency table are statistically independent. Thus,*

$$\text{rank} = \begin{cases} 2, & \text{dependent} \\ 1, & \text{statistical independent} \end{cases}$$

5.2 Independence of 3×3 Contingency Table

When the number of rows and columns are larger than 3, then the situation is a little changed. It is easy to see that the rank for statistical independence of a $m \times n$ contingency table is equal 1.0 as shown in Theorem 3. Also, when the rank is equal to $\min(m, n)$, two attributes are dependent.

Then, what kind of structure will a contingency matrix have when the rank is larger than 1,0 and smaller than $\min(m, n) - 1$? For illustration, let us consider the following 3times3 contingency table.

Example. Let us consider the following corresponding matrix:

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}.$$

The determinant of A is:

$$\begin{aligned} \det(A) &= 1 \times (-1)^{1+1} \det \begin{pmatrix} 5 & 6 \\ 8 & 9 \end{pmatrix} \\ &\quad + 2 \times (-1)^{1+2} \det \begin{pmatrix} 4 & 6 \\ 7 & 9 \end{pmatrix} \\ &\quad + 3 \times (-1)^{1+3} \det \begin{pmatrix} 4 & 5 \\ 7 & 8 \end{pmatrix} \\ &= 1 \times (-3) + 2 \times 6 + 3 \times (-3) = 0 \end{aligned}$$

Thus, the rank of A is smaller than 2. On the other hand, since $(123) \neq k(456)$ and $(123) \neq k(789)$, the rank of A is not equal to 1.0 Thus, the rank of A is equal to 2.0. Actually, one of three rows can be represented by the other two rows. For example,

$$(4\ 5\ 6) = \frac{1}{2}\{(1\ 2\ 3) + (7\ 8\ 9)\}.$$

Therefore, in this case, we can say that two of three pairs of one attribute are dependent to the other attribute, but one pair is statistically independent of the other attribute with respect to the linear combination of two pairs. It is easy to see that this case includes the cases when two pairs are statistically independent of the other attribute, but the table becomes statistically dependent with the other attribute.

In other words, the corresponding matrix is a mixture of statistical dependence and independence. We call this case *contextual independent*. From this illustration, the following theorem is obtained:

Theorem 6. *If the rank of the corresponding matrix of a 3×3 contingency table is 1, then two attributes in a given contingency table are statistically independent. Thus,*

$$\text{rank} = \begin{cases} 3, & \text{dependent} \\ 2, & \text{contextual independent} \\ 1, & \text{statistical independent} \end{cases}$$

It is easy to see that this discussion can be extended into $3 \times n$ contingency tables.

5.3 Independence of $m \times n$ Contingency Table

Finally, the relation between rank and independence in a multi-way contingency table is obtained from Theorem 3.

Theorem 7. *Let the corresponding matrix of a given contingency table be a $m \times n$ matrix. If the rank of the corresponding matrix is 1, then two attributes in a given contingency table are statistically independent. If the rank of the corresponding matrix is $\min(m, n)$, then two attributes in a given contingency table are dependent. Otherwise, two attributes are contextual dependent, which means that several conditional probabilities can be represented by a linear combination of conditional probabilities. Thus,*

$$\text{rank} = \begin{cases} \min(m, n) & \text{dependent} \\ 2, \dots, \min(m, n) - 1 & \text{contextual independent} \\ 1 & \text{statistical independent} \end{cases}$$

6 Conclusion

In this paper, a contingency table is interpreted from the viewpoint of granular computing and statistical independence. From the definition of statistical independence, statistical independence in a contingency table will hold when the equations of collinearity (Equation 6) are satisfied. In other words, statistical independence can be viewed as linear dependence. Then, the correspondence between contingency table and matrix, gives the theorem where the rank of the contingency matrix of a given contingency table is equal to 1 if two attributes are statistical independent. That is, all the rows of contingency table can be described by one row with the coefficient given by a marginal distribution. If the rank is maximum, then two attributes are dependent. Otherwise, some probabilistic structure can be found within attribute-value pairs in a given attribute, which we call contextual independence. Thus, matrix algebra is a key point of the analysis of a contingency table and the degree of independence, rank plays a very important role in extracting a probabilistic model.

References

1. Butz, C.J. Exploiting contextual independencies in web search and user profiling, *Proceedings of World Congress on Computational Intelligence (WCCI'2002)*, CD-ROM, 2002.
2. Coxeter, H.S.M. *Projective Geometry, 2nd Edition*, Springer, New York, 1987.
3. Polkowski, L. and Skowron, A. (Eds.) *Rough Sets and Knowledge Discovery 1*, Physica Verlag, Heidelberg, 1998.
4. Polkowski, L. and Skowron, A. (Eds.) *Rough Sets and Knowledge Discovery 2*, Physica Verlag, Heidelberg, 1998.
5. Pawlak, Z., *Rough Sets*. Kluwer Academic Publishers, Dordrecht, 1991.
6. Rao, C.R. *Linear Statistical Inference and Its Applications, 2nd Edition*, John Wiley & Sons, New York, 1973.
7. Skowron, A. and Grzymala-Busse, J. From rough set theory to evidence theory. In: Yager, R., Fedrizzi, M. and Kacprzyk, J. (eds.) *Advances in the Dempster-Shafer Theory of Evidence*, pp.193-236, John Wiley & Sons, New York, 1994.
8. Tsumoto S and Tanaka H: Automated Discovery of Medical Expert System Rules from Clinical Databases based on Rough Sets. In: *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining 96*, AAAI Press, Palo Alto CA, pp.63-69, 1996.
9. Tsumoto, S. Knowledge discovery in clinical databases and evaluation of discovered knowledge in outpatient clinic. *Information Sciences*, **124**, 125-137, 2000.
10. Yao, Y.Y. and Wong, S.K.M., A decision theoretic framework for approximating concepts, *International Journal of Man-machine Studies*, **37**, 793-809, 1992.
11. Yao, Y.Y. and Zhong, N., An analysis of quantitative measures associated with rules, N. Zhong and L. Zhou (Eds.), *Methodologies for Knowledge Discovery and Data Mining, Proceedings of the Third Pacific-Asia Conference on Knowledge Discovery and Data Mining*, LNAI **1574**, Springer, Berlin, pp. 479-488, 1999.
12. Ziarko, W., Variable Precision Rough Set Model. *Journal of Computer and System Sciences*, **46**, 39-59, 1993.