

Information theoretic models for dependence analysis and missing data estimation

D. S. Hooda, Permil Kumar

Department of Mathematics, Jaypee University of Engineering and Technology, A.B. Road, Raghogarh, Distt.Guna-473226 (M.P.) India
 Department of Statistics, University of Jammu, Jammu-(India)

Email address:

ds_hooda@rediffmail.com (D. S. Hooda), Parmil@yahoo.com (P. Kumar)

To cite this article:

D.S.Hooda, Permil Kumar. Information Theoretic Models for Dependence Analysis And missing Data Estimation, *American Journal of Theoretical and Applied Statistics*. Vol. 2, No. 2, 2013, pp. 15-20. doi: 10.11648/j.ajtas.20130202.12

Abstract: In the present communication information theoretic dependence measure has been defined using maximum entropy principle, which measures amount of dependence among the attributes in a contingency table. A relation between information theoretic measure of dependence and Chi-square statistic has been discussed. A generalization of this information theoretic dependence measure has been also studied. In the end Yate’s method and maximum entropy estimation of missing data in design of experiment have been described and illustrated by considering practical problems with empirical data.

Keywords: Maximum Entropy Principle, Contingency Table, Chi-Square Statistics, Lagrange’s Multipliers And Dependence Measure

1. Introduction

The frequencies of data based on counting of objects or units are categorized in a classified table known as contingency table. It can be defined as a rectangular array of order (m×n) having mn cells, where m and n are the number of rows and columns, which are equal to the number of categories of two attributes. In matrix notation, we have the following Contingency Table:

Table

| | Attribute B | | | | Total |
|----------------|-----------------|-----------------|-----------------|-----------------|----------------|
| | B ₁ | B ₂ | B _i | B _n | |
| A ₁ | O ₁₁ | O ₁₂ | O _{2i} | O _{1n} | r ₁ |
| A ₂ | O ₂₁ | O ₂₂ | O _{2i} | O _{2n} | r ₂ |
| B _i | O _{i1} | O _{i2} | O _{ii} | O _{in} | r _i |
| B _m | O _{m1} | O _{m2} | O _{mi} | O _{mn} | r _m |
| total | c ₁ | c ₂ | c _i | c _n | T |

In above table, O_{ij}, the (i,j)th cell represents the frequency of characteristics of A_i and B_j, and r_i and c_j are the marginal row and column sum totals.

The null hypothesis H₀ of independence of attributes against H₁ that attributes are dependent can be tested by

chi-square test statistic i.e.

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{(o_{ij} - e_{ij})^2}{e_{ij}} \tag{1}$$

where e_{ij} is the expected frequency corresponding to (i,j)th cell having observed frequency O_{ij}. Under the null hypothesis of independence

$$e_{ij} = \frac{r_i c_j}{T}$$

where

$$T = \sum_{i=1}^m r_i = \sum_{j=1}^n c_j \tag{2}$$

A decision about H₀ is made by comparing the value of calculated chi-square i.e. (1) with the tabulated value of chi-square for (m-1) (n-1) degrees of freedom at α % level of significance.

In section 2, an information theoretic dependence measure has been derived by using maximum entropy principle, which measures amount of dependence among the attributes or attributes of the contingency table. A relationship between information theoretic measure of dependence

and χ^2 statistic has been discussed in section 3. In section 4, a generalized information theoretic dependence measure has been studied. Yate's method and maximum entropy estimation of missing data in design of experiments have been described in section 5.

2. Information Theoretic Dependence Measure

Contingency table paves the way for the analysis of categorical data in physical and social sciences. But, many times only row and column totals are provided. In such cases we make use of Maximum Entropy Principle (MEP) which gives the same estimate as given by the hypothesis of independence. Soofi and Gokhale (1997) presented an information theoretic formulation, which gave an insight into the degree of dependence among the factors of contingency table. Let O_{ij} be the observed frequency in i th row and j th column of the $m \times n$ contingency table. Let r_1, r_2, \dots, r_m and c_1, c_2, \dots, c_n be the row and column sums or row and column marginal totals such that

$$\sum_{j=1}^n O_{ij} = r_i, i = 1, 2, \dots, m \quad (3)$$

$$\sum_{i=1}^m O_{ij} = c_j, j = 1, 2, \dots, n \quad (4)$$

And

$$\sum_{i=1}^m \sum_{j=1}^n O_{ij} = \sum_{i=1}^m r_i = \sum_{j=1}^n c_j = T \quad (5)$$

In case only the marginal totals information is provided to us, then the cell frequencies have to be estimated. We can fill only $(m-1)(n-1)$ cell frequencies arbitrarily and determine the $mn - (m-1)(n-1) = m + n - 1$ cell frequencies subject to (3), (4) and (5). Thus, there can be infinite number of sets that will be consistent with the given row and column totals. Out of these, we choose one which has maximum entropy as the above described situation demands the use of MEP i.e. we should choose the cell values so as to maximize

$$S = - \sum_{i=1}^m \sum_{j=1}^n \frac{x_{ij}}{T} \log \frac{x_{ij}}{T} \quad (6)$$

Subject to

$$\sum_{j=1}^n x_{ij} = r_i, i = 1, 2, \dots, m \quad (7)$$

$$\sum_{i=1}^m x_{ij} = c_j, j = 1, 2, \dots, n \quad (8)$$

And

$$\sum_{i=1}^m \sum_{j=1}^n x_{ij} = T \quad (9)$$

Using Lagrange's method of multipliers, we have

$$L = - \sum_{i=1}^m \sum_{j=1}^n \frac{x_{ij}}{T} \log \frac{x_{ij}}{T} - \lambda_0 \left(\sum_{j=1}^n x_{ij} - r_i \right) - \lambda_1 \left(\sum_{i=1}^m x_{ij} - c_j \right) - \lambda_2 \left(\sum_{i=1}^m \sum_{j=1}^n x_{ij} - T \right), \quad (10)$$

where λ_0, λ_1 and λ_2 are undetermined Lagrange's multipliers.

Differentiating (10), we get

$$\frac{\partial L}{\partial x_{ij}} = - \frac{1}{T} \log \frac{x_{ij}}{T} - \frac{1}{T} - \lambda_0 - \lambda_1 - \lambda_2 \quad (11)$$

Equating (11) equal to zero, we get

$$\hat{x}_{ij} = T e^{T(-\lambda_0 - \lambda_1 - \lambda_2 - \frac{1}{T})} \quad (12)$$

(7), (8) and (9) together with (12) respectively give

$$nT e^{T(-\lambda_0 - \lambda_1 - \lambda_2 - \frac{1}{T})} = r_i \quad (13)$$

$$mT e^{T(-\lambda_0 - \lambda_1 - \lambda_2 - \frac{1}{T})} = c_j \quad (14)$$

$$mnT e^{T(-\lambda_0 - \lambda_1 - \lambda_2 - \frac{1}{T})} = T \quad (15)$$

From (13) and (14), we get

$$mnT e^{T(-\lambda_0 - \lambda_1 - \lambda_2 - \frac{1}{T})} \cdot T e^{T(-\lambda_0 - \lambda_1 - \lambda_2 - \frac{1}{T})} = c_j r_i \quad (16)$$

Using (15) in (16), we have

$$T e^{T(-\lambda_0 - \lambda_1 - \lambda_2 - \frac{1}{T})} = \frac{r_i c_j}{T} \quad (17)$$

Equation (17) together with (12) gives

$$\hat{x}_{ij} = \frac{r_i c_j}{T}, \quad (18)$$

which is the maximum entropy estimate of the (i,j) th cell frequency. Let us denote this maximum entropy estimate by e_{ij} . Then, maximum entropy is

$$\begin{aligned} S_{\max} &= - \sum \sum \frac{\hat{x}_{ij}}{T} \log \frac{\hat{x}_{ij}}{T} \\ &= - \sum \sum \frac{e_{ij}}{T} \log \frac{e_{ij}}{T} \\ &= - \sum \sum \frac{r_i c_j}{T^2} \log \frac{r_i c_j}{T^2} \\ &= - \sum_i \frac{r_i}{T} \log \frac{r_i}{T} - \sum_j \frac{c_j}{T} \log \frac{c_j}{T} \\ &= S_r + S_c, \end{aligned} \quad (19)$$

where S_r and S_c are the entropies of row and column to-

tals respectively.

Let O_{ij} 's be the observed cell frequencies, then the entropy is given by

$$S = - \sum_i \sum_j \frac{O_{ij}}{T} \log \frac{O_{ij}}{T} \tag{20}$$

Since $S_{\max} \geq S$, therefore,

$$S_r + S_c \geq S$$

It implies

$$D = S_r + S_c - S \geq 0, \tag{21}$$

where D is the difference between S_{\max} and S and is called information theoretic measure of dependence [refer to Watanabe (1969) and Kapur and Kesavan (1992)].

Actually, $D > 0$, measures the information contained in contingency table in addition to information given by row and column sum totals.

From (21), we know

$$\begin{aligned} D &= - \sum_i \frac{r_i}{T} \log \frac{r_i}{T} - \sum_j \frac{c_j}{T} \log \frac{c_j}{T} + \\ &\quad + \sum_i \sum_j \frac{O_{ij}}{T} \log \frac{O_{ij}}{T} \\ &= - \sum_i \sum_j \frac{O_{ij}}{T} \log \frac{r_i}{T} \frac{c_j}{T} + \sum_i \sum_j \frac{O_{ij}}{T} \log \frac{O_{ij}}{T} \\ &= \sum_i \sum_j \frac{O_{ij}}{T} \log \frac{O_{ij}}{e_{ij}} \end{aligned} \tag{22}$$

It vanishes if and only if

$$\frac{O_{ij}}{T} = \frac{r_i}{T} \frac{c_j}{T}$$

Or

$$O_{ij} = \frac{r_i c_j}{T} = e_{ij}$$

Thus, we can conclude that the maximum entropy estimates e_{ij} are equal to the estimates obtained from the hypothesis of independence. The amount of deficit or difference between O_{ij} and e_{ij} values is due to the association or dependence between the factors of contingency tables. So, D is an information theoretic measure, which can be used for measuring the dependence between the factors of contingency table.

3. Information Theoretic Dependence Measure and χ^2

In this section we study the relationship between information theoretic dependence measure D and χ^2 statistic.

Let

$$O_{ij} = e_{ij} + \epsilon_{ij} \tag{23}$$

where ϵ_{ij} is very small quantity,

Since

$$\sum_{i=1}^m \sum_{j=1}^n O_{ij} = \sum_{i=1}^m \sum_{j=1}^n e_{ij} = T \tag{24}$$

And

$$\sum_{i=1}^m \sum_{j=1}^n O_{ij} = \sum_{i=1}^m \sum_{j=1}^n e_{ij} + \sum_{i=1}^m \sum_{j=1}^n \epsilon_{ij}$$

therefore

$$\sum_{i=1}^m \sum_{j=1}^n \epsilon_{ij} = 0 \tag{25}$$

(22) together with (23) gives

$$\begin{aligned} D &= \sum_{i=1}^m \sum_{j=1}^n \frac{e_{ij} + \epsilon_{ij}}{T} \log \left(\frac{e_{ij} + \epsilon_{ij}}{e_{ij}} \right) \\ &= \sum_{i=1}^m \sum_{j=1}^n \frac{e_{ij} + \epsilon_{ij}}{T} \log \left(1 + \frac{\epsilon_{ij}}{e_{ij}} \right) \\ &= \sum_{i=1}^m \sum_{j=1}^n \frac{e_{ij} + \epsilon_{ij}}{T} \left[\frac{\epsilon_{ij}}{e_{ij}} - \frac{1}{2} \frac{\epsilon_{ij}^2}{e_{ij}^2} + \frac{1}{3} \frac{\epsilon_{ij}^3}{e_{ij}^3} + \dots \right] \\ &= \sum_{i=1}^m \sum_{j=1}^n \frac{\epsilon_{ij}}{T} + \sum_{i=1}^m \sum_{j=1}^n \frac{1}{2T} \frac{\epsilon_{ij}^2}{e_{ij}} + \dots \end{aligned}$$

Using (25) and neglecting higher order terms, we have

$$\begin{aligned} D &\sim \sum_{i=1}^m \sum_{j=1}^n \frac{1}{2T} \frac{\epsilon_{ij}^2}{e_{ij}} \\ &= \frac{1}{2T} \sum_{i=1}^m \sum_{j=1}^n \frac{(O_{ij} - e_{ij})^2}{e_{ij}} \\ &= \frac{1}{2T} \chi^2 \end{aligned} \tag{26}$$

Thus, D gives times chi-square information or χ^2 gives $2T$ times the additional information given by the observed frequencies over the information already provided by the row and column totals. It is worth mentioning that χ^2 test does not give us the degree of dependence, while it is provided by information theoretic measure of dependence D . Moreover, χ^2 statistic cannot be used for comparing the several tables. Many coefficients have been proposed to measure the association between row and column factors, viz., Yule's coefficient Q of association, Pearson's coefficients ϕ^2 of mean square contingency, etc.

Similarly, D can be normalized to make it applicable to compare the dependence of factors of several tables. The

normalized D is given by

$$D_N = e^{-D}$$

For $D = 0$, D_N takes value one and for very large D , D_N is 0. It implies $0 \leq D_N \leq 1$. It may be noted that when $D_N = 1$, attributes or factors are independent and for $D_N = 0$, factors are perfectly dependent.

4. Generalized Measure of Dependence

In the present section, we study a generalized measure of dependence in contingency table of which D given by (2.19) is a particular case.

We choose the cell frequency which maximizes Harvda and Charvat (1967) entropy of degree β given below:

$$S^\beta = \frac{1}{1-\beta} \left[\sum_{i=1}^m \sum_{j=1}^n \left(\frac{x_{ij}}{T} \right)^\beta - 1 \right], \beta(>0) \neq 1 \quad (27)$$

subject to

$$\sum_{j=1}^n x_{ij} = r_i \quad (28)$$

$$\sum_{i=1}^m x_{ij} = c_j \quad (29)$$

and

$$\sum_{i=1}^m \sum_{j=1}^n x_{ij} = T \quad (30)$$

We use Lagrange's method of multipliers to maximize (27) subject to constraints (28) to (30).

$$\begin{aligned} L = & \frac{1}{1-\beta} \left[\left(\frac{x_{ij}}{T} \right)^\beta - 1 \right] + \lambda_0 \left(\sum_{j=1}^n x_{ij} - r_i \right) \\ & + \lambda_1 \left(\sum_{i=1}^m x_{ij} - c_j \right) + \lambda_2 \left(\sum_{i=1}^m \sum_{j=1}^n x_{ij} - T \right), \quad (31) \end{aligned}$$

where λ_0 , λ_1 and λ_2 are undetermined Lagrange's multipliers. Now, differentiating (31) w.r.t. x_{ij} and equating it to zero, we get

$$\frac{\partial L}{\partial x_{ij}} = \frac{1}{1-\beta} \left[\frac{1}{T^\beta} \cdot \beta \cdot x_{ij}^{\beta-1} + \lambda_0 + \lambda_1 + \lambda_2 \right] = 0$$

It implies

$$\frac{\beta}{1-\beta} \frac{x_{ij}^{\beta-1}}{T^\beta} = -\lambda_0 - \lambda_1 - \lambda_2$$

$$x_{ij}^{\beta-1} = \frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2)$$

Or

$$\hat{x}_{ij} = \left[\frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2) \right]^{\frac{1}{\beta-1}} \quad (32)$$

where λ_0 , λ_1 and λ_2 to be determined using (28), (29) and (30). Equations (28), (29) and (30) together with (32) respectively give

$$n \left[\frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2) \right]^{\frac{1}{\beta-1}} = r_i \quad (33)$$

$$m \left[\frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2) \right]^{\frac{1}{\beta-1}} = c_j \quad (34)$$

and

$$mn \left[\frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2) \right]^{\frac{1}{\beta-1}} = T \quad (35)$$

From (33) and (34), we have

$$\begin{aligned} mn \left[\frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2) \right]^{\frac{1}{\beta-1}} \left[\frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2) \right]^{\frac{1}{\beta-1}} \\ = r_i c_j \quad (36) \end{aligned}$$

From (36) and (35) we get

$$\left[\frac{\beta-1}{\beta} T^\beta (\lambda_0 + \lambda_1 + \lambda_2) \right]^{\frac{1}{\beta-1}} = \frac{r_i c_j}{T} \quad (37)$$

On putting the value of (37) in (32), we have

$$\hat{x}_{ij} = \frac{r_i c_j}{T} \quad (38)$$

Thus, (38) is the maximum entropy estimate of the (i, j)th cell frequency denoted by e_{ij} . Hence

$$\begin{aligned} S^\beta_{\max} &= \frac{1}{1-\beta} \left[\sum_{i=1}^m \sum_{j=1}^n \left(\frac{e_{ij}}{T} \right)^\beta - 1 \right] \\ &= \frac{1}{1-\beta} \left[\sum_{i=1}^m \sum_{j=1}^n \left(\frac{r_i}{T} \frac{c_j}{T} \right)^\beta - 1 \right] \\ &= S_r^\beta \sum_{j=1}^n \left(\frac{c_j}{T} \right) + S_c^\beta \quad (39) \end{aligned}$$

where

$$S_r^\beta = \frac{1}{1-\beta} \left[\sum_{i=1}^m \left(\frac{r_i}{T} \right) - 1 \right] \quad (40)$$

is β degree entropy of row total and

$$S_c^\beta = \frac{1}{1-\beta} \left[\sum_{j=1}^n \left(\frac{c_j}{T} \right) - 1 \right] \quad (41)$$

is β degree entropy of column total

From (41), we have

$$\sum_{j=1}^n \left(\frac{c_j}{T}\right) = (1-\beta)S_c^\beta + 1$$

Putting this value in (39), we get

$$\begin{aligned} S_{\max}^\beta &= S_r^\beta (1 + (1-\beta)S_c^\beta) + S_c^\beta \\ &= S_r^\beta + (1-\beta)S_r^\beta S_c^\beta + S_c^\beta \\ &= S_r^\beta + S_c^\beta + (1-\beta)S_r^\beta S_c^\beta \end{aligned} \tag{42}$$

Since $S_{\max}^\beta \geq S^\beta$, therefore,

$$S_r^\beta + S_c^\beta + (1-\beta)S_r^\beta S_c^\beta \geq S^\beta$$

Hence

$$D_\beta = S_r^\beta + S_c^\beta + (1-\beta)S_r^\beta S_c^\beta - S^\beta \geq 0, \tag{43}$$

which is equal to 0 if $S_{\max}^\beta = S^\beta$. Thus, D_β is generalized information theoretic dependence measure between the factors of contingency tables and reduce to D given by (21) in case $\beta \rightarrow 1$.

5. Estimation of Missing Data in Design of Experiments

In field experiments we design the field plots. In case we find one or more observations missing due to natural calamity or destroyed by a pest or eaten by animals, it is cumbersome to estimate the missing value or values as in field trials it is practically impossible to repeat the experiment under identical conditions. So we have no option except to make best use of the data available. Yates (1933) suggested a method:

“Substitute x for the missing value and then choose x so as to minimize the error sum of squares”.

Actually, the substituted value does not recover the best information, however, it gives the best estimate according to a criterion based on the least square method.

For the randomized block experiment

$$x = \frac{pP + qQ - T}{(p-1)(q-1)} \tag{44}$$

where

p = Number of treatments,

q = Number of blocks

P = Total of all plots receiving the same treatment as the missing plot

Q = Total of all plots in the same block as the missing plot

T = Total of all plots

For the Latin Square Design, the corresponding formula is

$$x = \frac{p(P_r + P_c + P_t) - 2T}{(p-1)(q-1)} \tag{45}$$

where

P = Number of rows or columns of treatments

P_r = Total of row containing the missing plot

P_c = Total of column containing the missing plot

P_t = Total of treatment contained in the missing plot

T = Grand total

In case more than one plot yields are missing, we substitute the average yield

of available plots in all except one of these and substitute x in this plot. We estimate x by Yate’s method and use this value to estimate the yields of other plots one by one.

Next we discuss the maximum entropy method. If

x_1, x_2, \dots, x_n are known

yields and x is the missing yield. We obtain the maximum entropy estimate

for x by maximizing:

$$-\sum_{i=0}^n \frac{x_i}{T+x} \log \frac{x_i}{T+x} - \frac{x}{T+x} \log \frac{x}{T+x} \tag{46}$$

Thus we get

$$\hat{x} = \left[x_1^{x_1} x_2^{x_2} \dots x_n^{x_n} \right]^{\frac{1}{T}}, \tag{47}$$

where $T = \sum_{i=1}^n x_i$ and the value given by (47) is called

maximum entropy mean of x_1, x_2, \dots, x_n .

Similarly, if two values x and y are missing, x and y are determined from

$$\hat{x} = \left[x_1^{x_1} x_2^{x_2} \dots x_n^{x_n} \right]^{\frac{1}{T+y}}, \tag{48}$$

$$\hat{y} = \left[x_1^{x_1} x_2^{x_2} \dots x_n^{x_n} \right]^{\frac{1}{T+x}}, \tag{49}$$

The solution of (48) and (49) is

$$\hat{x} = \hat{y} = \left[x_1^{x_1} x_2^{x_2} \dots x_n^{x_n} \right]^{\frac{1}{T}} \tag{50}$$

Hence all the missing values have the same estimate and this does not change if the missing values are estimated one by one.

There are three following drawbacks of the estimate given by (47)

(i) \hat{x} is rather unnatural. In fact \hat{x} is always greater than Arithmetic mean of x_1, x_2, \dots, x_n .

(ii) If two values are missing, the maximum entropy estimated for each is the same as given by (50).

(iii) This is not very useful for estimating missing values in design of experiments.

The first drawback can be overcome by using generalized measure of entropy instead of Shannon entropy. If we

use Burg's [1] measure given by

$$B(P) = \sum_{i=1}^n \log p_i \quad (51)$$

Then we get the estimate

$$\hat{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \bar{x} \quad (52)$$

In fact we choose a value \hat{x} , which is as equal to x_1, x_2, \dots, x_n as possible and so we maximize a measure of equality. Since there are many measures of equality, therefore our estimate will also depend on the measure of equality we choose.

The second drawback can be understood by considering the fact that the information theoretic estimate for a missing value depends on:

- (a) The information available to us
- (b) The purpose for which missing value is to be used.

The third drawback, as according to the principle of maximum entropy, we should use all the information given to us and avoid scrupulously using any information not given to us. But in design of experiments, we are given information about the structure of the design which we are not using this knowledge in estimating the missing values. Consequently, the estimate is not accurate, however, method defined in section 2 be applied to estimate the missing val-

ue x_{ij} in contingency tables. Accordingly, the value x_{ij} is to be chosen to minimize the measure of dependence D.

References

- [1] Burg J.P.(1970). "The relationship between Maximum Entropy Spectra and Maximum Likelihood in Modern Spectra Analysis", ed D.G. Childers, pp130-131.
- [2] Harvda, J. and Charvat, F. (1967). Quantification method of classification processes concepts of structural - entropy, *Kybernetika*, 3: 30-35.
- [3] Kapur, J.N. and Kesavan, H.K. (1992). "Entropy optimization principles with applications." Academic press, San Diego. K. Elissa, "Title of paper if known," unpublished.
- [4] Soofi, E.S. and Gokhale, D.V. (1997). "Information theoretic methods for categorical data. *Advances in Econometrics*.", JAI Press, Greenwich.
- [5] Watanabe, S.(1969). "Knowing and Guessing". John Wiley, New York, 1969.
- [6] Watanabe, S. (1981). "Pattern recognition as a quest for minimum entropy.", *Pattern Recognition*. 13:381-387.
- [7] Yates, F. (1933). "The analysis of replicated experiments when the field experiments are incomplete.", *Emp. Journ. Exp. Agri.* 1, 129-142.