

## UNEQUAL SIZE, TWO-WAY ANALYSIS OF VARIANCE FOR CATEGORICAL DATA

BY HAN-YONG CHUNG

### 1. Introduction

The techniques about the analysis of variance for quantitative variables have been well-developed. But when the variable is categorical, we must switch to a completely different set of varied techniques. R. J. Light and B. H. Margolin [1] presented one kind of techniques for categorical data in their paper, where there are  $G$  unordered experimental groups and  $I$  unordered response categories.

Assume that there are  $n$  responses  $X_1, \dots, X_n$ . Each  $X_i$  is the name of one of  $I$  possible categories.

Definition; The variation for categorical responses  $X_1, \dots, X_n$  is

$$\frac{1}{2n} \sum_{j=1}^n \sum_{i=1}^n d_{ij}^2 = \frac{1}{2n} \sum_{j=1}^n \sum_{i=1}^n d_{ij}$$

where  $d_{ij}=1$  if  $X_i$  and  $X_j$  name different categories.

$=0$  if  $X_i$  and  $X_j$  name the same category.

For  $n$  responses, each in one and only one of  $I$  possible categories, the data can be summarized with a vector  $\Phi$  of category counts  $\Phi = (n_1, \dots, n_I)$ , where  $n_i$  is the number of responses in the  $i$ th category,  $i=1, \dots, I$ , so that  $\sum_{i=1}^I n_i = n$ . The variation of these responses is:

$$\frac{1}{2n} [\sum_{j \neq i} n_i n_j] = \frac{1}{2n} [n^2 - \sum_{i=1}^I n_i^2]$$

To further motivate this definition of variation, we need the following known lemmas [1]:

LEMMA 1. *The variation of  $n$  categorical responses is minimized if and only if they all belong to the same category.*

LEMMA 2. *The variation of  $n$  responses, where  $n=IS+L$ ,  $0 \leq L < I$ , is*

maximized for any vector  $\Phi$  of category counts such that  $L$  counts equal  $S+1$ , and  $I-L$  counts equal  $S$ .

This note is an extension of one of the techniques to a two-way table, where there are  $I$  unordered response categories,  $J$  unordered experimental levels crossed by another  $K$  unordered experimental levels, with unequal size of observations in each of  $JK$  cells. For terminology and notation, we follow [1].

## 2. The model and variation components

We construct the two-way table where there are  $I$  unordered response categories,  $J$  unordered experimental levels crossed by another  $K$  unordered experimental levels with an unequal size of observations in each  $JK$  cells. Each response is in one and only one of the  $I$  categories. Denote the number of responses in category  $i$ ,  $j$ th level (of the second index),  $k$ th level (of the third index) by  $n_{ijk}$ .

We assume that responses in different cells are stochastically independent, and that each cell's responses  $(n_{1jk}, n_{2jk}, \dots, n_{Ijk})$  obey a multinomial law:

$$Pr\{(n_{1jk}, \dots, n_{Ijk})\} = \binom{n_{\cdot jk}}{n_{1jk}, \dots, n_{Ijk}} \prod_{i=1}^I (p_{ijk})^{n_{ijk}}$$

where  $\sum_{i=1}^I p_{ijk} = 1$ ,  $p_{ijk} > 0$ ,  $i=1, \dots, I$ ,  $j=1, \dots, J$ , and  $k=1, \dots, K$ .

If we let

$$V = (n_{111}, n_{211}, \dots, n_{I11}, n_{121}, n_{221}, \dots, n_{I21}, \dots, n_{1J1}, n_{2J1}, \dots, n_{IJ1}, n_{112}, n_{212}, \dots, n_{I12}, \dots, n_{1JK}, n_{2JK}, \dots, n_{IJK}),$$

then

$$\begin{aligned} E(V) &= Y = (n_{\cdot 11}p_{111}, n_{\cdot 11}p_{211}, \dots, n_{\cdot 11}p_{I11}, \dots, n_{\cdot J1}p_{1J1}, \\ &\quad \dots, n_{\cdot J1}p_{2J1}, \dots, n_{\cdot J1}p_{IJ1}, n_{\cdot 12}p_{112}, n_{\cdot 12}p_{212}, \dots, n_{\cdot 12}p_{I12}, \\ &\quad \dots, n_{\cdot JK}p_{1JK}, n_{\cdot JK}p_{2JK}, \dots, n_{\cdot JK}p_{IJK})', \\ \text{Cov}(V) &= Z = Z_{11} \oplus Z_{21} \oplus \dots \oplus Z_{J1} \oplus Z_{12} \oplus Z_{22} \oplus \dots \oplus Z_{J2} \oplus \\ &\quad \dots \oplus Z_{1K} + \dots \oplus Z_{JK} \end{aligned}$$

where

$$Z_{jk} = n_{\cdot jk} \begin{pmatrix} p_{1jk}(1-p_{1jk}) & -p_{1jk}p_{2jk} & \dots & -p_{1jk}p_{Ijk} \\ & p_{2jk}(1-p_{2jk}) & \dots & -p_{2jk}p_{Ijk} \\ & & \ddots & \vdots \\ & & & \vdots \\ & & & \dots & p_{Ijk}(1-p_{Ijk}) \end{pmatrix}$$

and  $\oplus$  denotes the direct sum operation (see[2]).

With the two-way table introduced as our model we define the following variations:

The total variation in the response variable ( $TSS$ ) is

$$TSS = n/2 - \sum_{i=1}^I n_{i..}^2 / 2n;$$

the within-2nd index level variation ( $WSS_1$ ) is

$$WSS_1 = \sum_{j=1}^J (n_{.j.} / 2 - \sum_{i=1}^I n_{ij.}^2 / 2n_{.j.});$$

the between-2nd index level variation ( $BSS_1$ ) is

$$BSS_1 = TSS - WSS_1;$$

the within-3rd index level variation ( $WSS_2$ ) is

$$WSS_2 = \sum_{k=1}^K (n_{..k} / 2 - \sum_{i=1}^I n_{i.k}^2 / 2n_{..k});$$

the between-3rd index level variation ( $BSS_2$ ) is

$$BSS_2 = TSS - WSS_2;$$

the within-cell variation ( $WSS_3$ ) is

$$WSS_3 = \sum_{k=1}^K \sum_{j=1}^J (n_{.jk} / 2 - \sum_{i=1}^I n_{ijk}^2 / 2n_{.jk});$$

the between-cell variation ( $BSS_3$ ) is

$$BSS_3 = TSS - WSS_3;$$

where

$$\begin{aligned} n_{.jk} &= \sum_{i=1}^I n_{ijk}, & n_{i.k} &= \sum_{j=1}^J n_{ijk}, & n_{ij.} &= \sum_{k=1}^K n_{ijk}, \\ n_{i..} &= \sum_{j=1}^J \sum_{k=1}^K n_{ijk}, & n_{.j.} &= \sum_{i=1}^I \sum_{k=1}^K n_{ijk}, \\ n_{..k} &= \sum_{i=1}^I \sum_{j=1}^J n_{ijk}, & n &= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K n_{ijk} \end{aligned}$$

### 3. Definitions

DEFINITION 1. The interaction between the 2nd index level and the 3rd index level is defined as  $I = ESS_3 - ESS_1 - BSS_2$ .

DEFINITION 2.  $\mathcal{Q}$  is the space where  $I = 0$ .

DEFINITION 3.  $p_i$  is the probability of an element belonging to  $i$ th category.  $p_{ij.}$  is the probability of an element belonging to  $i$ th category and  $j$ th level, regardless of the 3rd index level.  $p_{i.k}$  is the probability of an element be-

longing to  $i$ th category and  $k$ th level, regardless of the 2nd index level.

DEFINITION 4. The hypothesis  $H_1$  is  $p_{ij} = p_i$  for all  $j$ . The hypothesis  $H_2$  is  $p_{i.k} = p_i$  for all  $k$ . The hypothesis  $H_3$  is  $p_{ijk} = p_i$  for all  $j$  and  $k$ .

#### 4. Testing of the hypothesis

THEOREM 4-1. (a) Under the hypothesis  $H_1$ ,

$$(n-1)(J-1)BSS_1/TSS$$

is asymptotically approximated as  $\chi^2_{(I-1)(J-1)}$ .

(b) Under the hypothesis  $H_2$ ,

$$(n-1)(I-1)BSS_2/TSS$$

is asymptotically approximated as  $\chi^2_{(I-1)(K-1)}$ .

*Proof.* The above facts can be proved as in the case of one-way table (see [1]). To prove (a), since there are  $I$  categories and  $J$  levels the degree of freedom is  $(I-1)(J-1)$ . (b) can be proved in the similar way.

THEOREM 4-2. With large  $n_{.jk} = n_{.j}n_{..k}/n$  for all  $j, k$ ,  $BSS_1$  and  $BSS_2$  are asymptotically independent under the hypothesis  $H_3$ .

*Proof.* With large  $n_{.jk}$ ,  $V$  is asymptotically multivariate normal, i.e.,  $V \sim N(Y, Z)$ . Under the hypothesis  $H_3$ ,  $Z$  can be reduced as

$$Z = Z_{11} \oplus Z_{21} \oplus \cdots \oplus Z_{jk} \oplus \cdots \oplus Z_{JK},$$

where

$$Z_{jk} = n_{.jk} \begin{pmatrix} p_1(1-p_1) & -p_1p_2 & \cdots & -p_1p_I \\ & p_2(1-p_2) & \cdots & -p_2p_I \\ & & \vdots & \vdots \\ & & & \vdots \\ & & \cdots & p_I(1-p_I) \end{pmatrix}$$

Let

$$T = -(U_{JK} \otimes I_I) / 2n, \quad A = Z_K \otimes I_{IJ}, \quad A' = X_K \otimes I_{IJ},$$

$$W_1 = -\frac{1}{2} \left( \frac{1}{n_{.1}} I_I \oplus \frac{1}{n_{.2}} I_I \oplus \cdots \oplus \frac{1}{n_{.j}} I_I \right),$$

$$B = I_K \otimes (Y_J \otimes I_I), \quad B' = I_K \otimes (Y_J \otimes I_I),$$

and

$$W_2 = -\frac{1}{2} \left( \frac{1}{n_{..1}} I_I \oplus \frac{1}{n_{..2}} I_I \oplus \cdots \oplus \frac{1}{n_{..K}} I_I \right),$$

where  $U_r$  is a  $r \times r$  matrix of ones,  $I_r$  is a  $r \times r$  identity matrix,  $X_r$  is a  $1 \times r$  matrix of ones, and  $Y_r$  is a  $r \times 1$  matrix of ones.

Then

$$\begin{aligned} TSS &= \frac{n}{2} + V' T V, & WSS_1 &= \frac{n}{2} + A' V W_1 A' V, \\ WSS_2 &= \frac{n}{2} + V' B W_2 B' V, & BSS_1 &= V' (T - A W_1 A') V, \\ BSS_2 &= V' (T - B W_2 B') V. \end{aligned}$$

Now to prove that  $BSS_1$  and  $BSS_2$  are independent, it suffices to show that:

$$(T - A W_1 A') Z (T - B W_2 B') = 0$$

(see[3]).

$$\begin{aligned} A W_1 A' &= -\frac{1}{2n} [U_K \otimes \left( \frac{n}{n_{..1}} I_I \oplus \frac{n}{n_{..2}} I_I \oplus \cdots \oplus \frac{n}{n_{..J}} I_I \right)], \\ B W_2 B' &= (U_J \otimes \frac{1}{n_{..1}} I_I) \oplus (U_J \otimes \frac{1}{n_{..2}} I_I) \oplus \cdots \oplus (U_J \otimes \frac{1}{n_{..K}} I_I), \\ (T - A W_1 A') Z (T - B W_2 B') &= Y_K \otimes (e_{XY}), \\ X &= 1, 2, \dots, IJ, \text{ and } Y = 1, 2, \dots, IJK. \end{aligned}$$

Here

$$e_{XY} = \begin{cases} \hat{p}_{s'}(1 - \hat{p}_{t'}) \left( \frac{n^2 n_{..st}}{n_{..s} n_{..t}} - n \right) & \text{if } s' = t', \\ \hat{p}_{s'} \hat{p}_{t'} \left( \frac{n^2 n_{..st}}{n_{..s} n_{..t}} - n \right) & \text{if } s' \neq t', \end{cases}$$

where

$$\begin{aligned} s' &= X - I \left[ \frac{X-1}{I} \right], & t' &= Y - I \left[ \frac{Y-1}{I} \right], \\ s &= \left[ \frac{X-1}{I} \right] + 1 - J \left[ \frac{X-1}{J} \right], & t &= \left[ \frac{Y-1}{IJ} \right] + 1. \end{aligned}$$

Since  $n_{..jk} = n_{..j} n_{..k} / n$  for all  $j, k$ ,  $\frac{n^2 n_{..st}}{n_{..s} n_{..t}} = n$ , i.e.,  $e_{XY} = 0$  for all  $X, Y$ .

Therefore,  $BSS_1$  and  $BSS_2$  are asymptotically independent.

THEOREM 4-3. *With large  $n_{jk}$ , in the space  $\Omega$ , and under the hypotheses  $H_1$ ,  $H_2$ , and  $H_3$ ,*

$$(n-1)(I-1)BSS_3/TSS$$

*is approximated as  $\chi^2_{(I-1)(J-1+K-1)}$ .*

*Proof.* If  $I=0$ , then  $BSS_3=BSS_1+BSS_2$ . Hence,

$$\frac{(n-1)(I-1)BSS_3}{TSS} = \frac{(n-1)(I-1)BSS_1 + (n-1)(I-1)BSS_2}{TSS}$$

With large  $n_{jk}$  and under the hypotheses  $H_1$  and  $H_2$ , the distributions of  $\frac{(n-1)(I-1)BSS_1}{TSS}$  and  $\frac{(n-1)(I-1)BSS_2}{TSS}$  are approximated as  $\chi^2_{(I-1)(J-1)}$  and  $\chi^2_{(I-1)(K-1)}$  respectively. With large  $n_{jk}$  and under the hypothesis  $H_3$ ,  $BSS_1$  and  $BSS_2$  are asymptotically independent. So  $(n-1)(I-1)BSS_3/TSS$  is approximated as  $\chi^2_{(I-1)(J-1+K-1)}$ .

### References

- [1] Richard J. Light and Barry H. Margolin, *An analysis of Variance for Categorical Data*, *Journal of the American Statistical Association*, Vol. 66, No. 335 (1971), 534-544.
- [2] Graybill, F.A., *Introduction to Matrices with Applications in Statistics*, Wadsworth Publishing Co. Inc., 1969.
- [3] Graybill, F.A. *An Introduction to Linear Statistical Models*, New York; McGraw-Hill Book Co., 1961.

Seoul National University