



1

22.25

Walkson al

MICROCOPY RESOLUTION TEST CHART NATIONAL BUREAU OF STANDARDS-1963-A AD-A164 638

MEASURES OF IMBALANCE FOR

UNBALANCED MODELS

Вy

Andre' I. Khuri

Department of Statistics, University of Florida Gainesville, FL. 32611

> Technical Report Number 255 January, 1986

UTIC FILE COPY

PREPARED UNDER GRANT NO. NOOO14-86-K-0059 FROM THE OFFICE OF NAVAL RESEARCH. ANDRE' I. KHURI AND RAMON C. LITTELL, PRINCIPAL INVESTIGATORS.

for public district



<mark>მ</mark> პ

2

24

187

SE	CUBITY	CL ASSIEIC	ATION OF	THIS	PAGE	(When L	ata Enlara	đ

.

REP	ORT DOCUMENTATI	ON PAGE	BEFORE COMPLETING FORM
REPORT NUMBER		2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
255			
200 TITLE (and Subtitie)			5 TYPE OF REPORT & PERIOD COVERED
MEACUDES OF	IMPALANCE FOR UNB	ALANCED MODELS	S THE OF HER ON THE PERIOD COVERED
MEASURES OF	Indalance For onde		
			6. PERFORMING ORG. REPORT NUMBER
AUTHOR(#)			8 CONTRACT OR GRANT NUMBER(.)
Andre' I. Kh	uri		N00014-86-K-0059
			(R&T 411455201)
Department of	Statistics	(6.58	AREA & WORK UNIT NUMBERS
Nuclear Scien	ces Center. Unive	rsitv of Florida	
Gainesville.	FL 32611		
L CONTROLLING OFF			12 REPORT DATE
Office of Nav	al Research		January 1986
Mathematical	Sciences Division	(Code 411)	13. NUMBER OF PAGES
Arlington, VA	22217-5000		21
4 MONITORING AGEN	CY NAME & ADDRESS(11 dil	ferent from Controlling Office)	15. SECURITY CLASS. (of this report)
			UNCLASS FIFT
			UNULASSIFIED
			154. DECLASSIFICATION DOWNGRADING SCHEDULE
DIST DIRITION STA	TENENT Stible Report		L
APPROVED FOR P	UBLIC RELEASE: D	ISTRIBUTION UNI, IMIT	ED m Report)
APPROVED FOR P	UBLIC RELEASE: D	ISTRIBUTION UNLIMIT	ED m Report)
APPROVED FOR P	UBLIC RELEASE: D	ISTRIBUTION UNLIMIT	ED m Report)
APPROVED FOR P	UBLIC RELEASE: D	ISTRIBUTION UNLIMIT	ED m Roport)
APPROVED FOR P	UBLIC RELEASE: D	ISTRIBUTION UNLIMIT	ED m Report)
APPROVED FOR P DISTRIBUTION STA SUPPLEMENTARY ( Unbalanced da Loglinear mod	UBLIC RELEASE: D TEMENT of the obstract onto NOTES ta, Nested models els, Chi-squared s	ISTRIBUTION UNLIMIT ered in Block 20, 11 different fro ry and identify by block number) , Cross-classificat statistic.	ED m Report) ion models,
APPROVED FOR P . DISTRIBUTION STA . SUPPLEMENTARY ( Unbalanced da Loglinear mode	UBLIC RELEASE: D TEMENT of the obstract only NOTES ta, Nested models els, Chi-squared s	ISTRIBUTION UNLIMIT ered in Block 20, 11 different fro ry end identify by block number) , Cross-classificat statistic.	ED m Report) ion models,
APPROVED FOR P DISTRIBUTION STA SUPPLEMENTARY I Unbalanced da Loglinear mode	UBLIC RELEASE: D	ISTRIBUTION UNLIMIT ered in Block 20, 11 different fro ry and identify by block number) , Cross-classificat statistic.	ED m Report) ion models,
APPROVED FOR P DISTRIBUTION STA SUPPLEMENTARY ( Unbalanced da Loglinear mod ABSTRACT (Continu Loglinear mod data set. T	UBLIC RELEASE: D TEMENT of the obstract onto NOTES NOTES els, Chi-squared s dels are used to m he proposed proced	ISTRIBUTION UNLIMIT ored in Block 20, 11 different fro ry end identify by block number) , Cross-classificat statistic. y end identify by block number) neasure the degree of lure can also be use	ED m Report) ion models, of imbalance of an unbalanced ed to measure departures
APPROVED FOR P DISTRIBUTION STA Supplementary ( Unbalanced da Loglinear mode ABSTRACT (Continu Loglinear mode data set. The from certain frequencies,	UBLIC RELEASE: D	ISTRIBUTION UNLIMIT ored in Block 20, 11 different fro ry end identify by block number) , Cross-classificat statistic. y end identify by block number) measure the degree of lure can also be use , such as proportion and last-stage unit	ED m Report) ion models, of imbalance of an unbalanced ed to measure departures hality of subclass formity.
APPROVED FOR P 2. DISTRIBUTION STA 3. SUPPLEMENTARY ( Unbalanced da Loglinear mode 4. Loglinear mode data set. The from certain frequencies,	UBLIC RELEASE: D	ISTRIBUTION UNLIMIT ored in Block 20, 11 different fro ry end identify by block number) , Cross-classificat statistic. y end identify by block number) measure the degree of lure can also be use , such as proportion and last-stage unit	ED m Report) ion models, of imbalance of an unbalanced ed to measure departures hality of subclass formity.
APPROVED FOR P 2. DISTRIBUTION STA 5. SUPPLEMENTARY ( Continuent of the set of the se	UBLIC RELEASE: D TEMENT of the obstract entry NOTES els, Chi-squared s els, Chi-squared s tapposed proced types of balance, partial balance,	ISTRIBUTION UNLIMIT or end identify by block number) , Cross-classificat statistic. y end identify by block number) measure the degree of lure can also be use , such as proportion and last-stage unit	ED m Report) ion models, of imbalance of an unbalanced ed to measure departures hality of subclass formity.
APPROVED FOR P . DISTRIBUTION STA . SUPPLEMENTARY ( Unbalanced da Loglinear mode data set. The from certain frequencies, D (JAN 73 1473	UBLIC RELEASE: D TEMENT of the obstract and NOTES on reverse aide if necessaring ta, Nested models; els, Chi-squared s els, Chi-squared s dels are used to m he proposed proced types of balance, partial balance, EDITION OF I NOV 65 IS OF	ISTRIBUTION UNLIMIT ored in Block 20, 11 different fro ry end identify by block number) , Cross-classificat statistic. y end identify by block number) neasure the degree of lure can also be use , such as proportion and last-stage unit and last-stage unit SSOLETE	ED m Report) ion models, of imbalance of an unbalanced ad to measure departures nality of subclass formity. <u>ASSIFIED</u> SSIFICATION OF THIS PAGE (When Data Enter SSIFICATION OF THIS PAGE (W

Measures of Imbalance For Unbalanced Models A. I. Khuri The University of Florida

Accession For
MTIS AT 21
DTIC T
Unampert á 🗌
Just: 1. 1. 1. 1
By
District ory
Avalletter Cres
petralia e in con
Dist Special



Technical Report Number 255

Department of Statistics University of Florida Gainesville, Florida 32611

January 1986

Measures of Imbalance For Unbalanced Models

A.I. Khuri

The University of Florida

#### Abstract

In this paper we present a procedure to measure the degree of imbalance of an unbalanced data set. The procedure is based on choosing an appropriate loglinear model for the subclass frequencies of the data. A measure of imbalance is then introduced as some function of the chi-squared statistic used in the goodness-of-fit test for the loglinear model. The proposed procedure can also be used to measure departures from certain types of balance, such as proportionality of subclass frequencies, partial balance, and last-stage uniformity.

Key words: Unbalanced data, nested models, cross-classification models, loglinear models, chi-squared statistic.

#### 1. Introduction

It is known that in a balanced data situation, parameter estimators and test statistics pertaining to the effects in the associated model have certain optimal properties. These properties, however, cannot be maintained once the data set becomes unbalanced. In this case, the statistical properties of the aforementioned estimators and test statistics will, to a large extent, depend on the pattern of the data subclass frequencies. Severe imbalance in the data can have adverse effects on the analysis, especially if that analysis is an adaptation of procedures pertaining to balanced data (see, for example,

-1 -

Cummings and Gaylor 1974).

Ahrens and Pincus (1981) presented two measures of imbalance for the oneway classification model. These measures were utilized to assess the efficiency of an associated unbalanced design as compared to a balanced design with the same number of observations. Other authors have alluded to the need to measure data imbalance; they include Hess (1979, p. 646) and Tietjen (1974, p. 576).

The purpose of this paper is to present a general procedure to measure imbalance of a data set for a given unbalanced model. It is shown that one of the two measures introduced by Ahrens and Pincus (1981) can be derived as a special case using this procedure. The proposed procedure can also be utilized to measure departures from certain types of balance other than complete balance where frequencies are equal within all the subclasses. These include partial balance, last-stage uniformity, and the case of proportional subclass frequencies in cross-classification models.

## 2. A General Procedure to Measure Imbalance

A measure of imbalance, denoted by  $\phi(D)$ , is a function of the subclass frequencies which are determined by the design D used in the experiment. This function takes values inside the closed interval [0,1]. Small values of  $\phi(D)$ indicate severe imbalance, whereas "near balance" cases are characterized by large values of  $\phi(D)$ . The data set is balanced if and only if  $\phi(D)=1$ . Furthermore, this function must remain invariant under any partial or complete replication of the design (see Ahrens and Pincus 1981).

The development of the function  $\phi(D)$  is based on the use of loglinear models. Several unbalanced models will be considered to illustrate the application of this procedure. These models include the one-way classification model, the two-way classification model, the three-way classification model.

-2-

the two-fold nested model, the three-fold nested model, and a model with a mixture of cross-classified and nested effects.

# 2.1 The One-Way Classification Model

Consider the one-way model

$$y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$
 (2.1)

 $(i = 1, 2, ..., a; j = 1, 2, ..., n_i)$ , where  $\mu$  is a fixed unknown parameter,  $\alpha_i$  is either a fixed parameter or a random variable, and  $\varepsilon_{ij}$  is a random error. Here D is the design D =  $\{n_1, n_2, ..., n_a\}$ .

We shall consider that the  $n_i$ 's have a multinomial distribution such that  $n_i$  has the binomial distribution  $B(n_i, \Pi_i)$ , where  $n_i = \sum_{i=1}^{a} n_i$  and  $\Pi_i$  is the probability of belonging to level i (i = 1,2,...,a). Hence,  $m_i = E(n_i) = n_i \Pi_i$  (i = 1,2,...,a). The  $m_i$ 's will be referred to as expected frequencies. On a logarithmic scale, the expected frequencies can be represented by the loglinear model

$$\log m_{i} = \bar{\mu} + \bar{\alpha}_{i}, \quad i = 1, 2, ..., a,$$
 (2.2)

where

$$\overline{\mu} = \log n_{i} + \frac{1}{a} \sum_{i=1}^{a} \log \Pi_{i}$$

$$\overline{\alpha}_{i} = \log \Pi_{i} - \frac{1}{a} \sum_{i=1}^{a} \log \Pi_{i}, \quad i = 1, 2, \dots, a.$$

We note that  $\sum_{i=1}^{a} \bar{\alpha}_{i} = 0$  and that model (2.2) is of the same form as model (2.1), except for the error term.

Let  $m_i$  denote the maximum likelihood estimate of  $m_i$  (i = 1,2,...,a). Under complete balance,  $\pi_i = 1/a$  for all i, hence  $m_i = n_i/a = n_i$ . Using Pearson's approximate chi-squared statistic for testing the hypothesis  $H_0:\pi_i = 1/a$  for all i we obtain

$$x^{2} = \sum_{i=1}^{a} (n_{i} - \bar{n}_{i})^{2} / \bar{n}_{i},$$

which under  $H_0$  has an asymptotic chi-squared distribution with  $\theta$  degrees of freedom, where, in general,  $\theta$  is the difference between the number of independent  $\Pi_i$ 's under  $H_a$  and under  $H_0$ , respectively. In this case  $\theta = a-1$ . We define our measure of imbalance as

$$\phi(D) = \frac{1}{1+c^2},$$
 (2.3)

where  $c^2 = X^2/n_{..}$  We note that  $0 \le \phi(D) \le 1$  and the division of  $X^2$  by n. causes the measure to be invariant to any replication of the design as required. Furthermore,  $\phi(D) = 1$  if and only if the n's are equal. We also note that  $\phi(D)$  is identical to the measure v(D) given by Ahrens and Pincus (1981).

## 2.2 The Two-Way Classification Model

Consider the model

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}, \qquad (2.4)$$

(i = 1,2,...,a; j = 1,2,...,b; k = 1,2,...,n<sub>ij</sub>), where  $\mu$  is a fixed unknown parameter;  $\alpha_i$  and  $\beta_j$  can be either fixed or random. In this case the design D is D = {n<sub>11</sub>,n<sub>12</sub>,...,n<sub>ab</sub>}. The n<sub>ij</sub>'s are considered to have the multinomial distribution and each n<sub>ij</sub> has the binomial distribution B(n<sub>..</sub>,  $\Pi_{ij}$ ), where n<sub>..</sub> =  $\Sigma_{i,j}n_{ij}$  and  $\Pi_{ij}$  is the probability of belonging to the (i,j)<sup>th</sup> cell. Hence, E(n<sub>ij</sub>) = m<sub>ij</sub> = n<sub>..</sub> $\Pi_{ij}$ . The corresponding loglinear model is

$$\log \mathfrak{m}_{ij} = \overline{\mu} + \overline{\alpha}_{i} + \overline{\beta}_{j} + (\overline{\alpha \overline{\beta}})_{ij}, \qquad (2.5)$$

where in this case

$$\overline{u} = \frac{1}{ab} \sum_{i,j} \log m_{ij}$$

$$\overline{\alpha}_{i} = \frac{1}{b} \sum_{j} \log m_{ij} - \overline{\mu}$$

$$\overline{\beta}_{j} = \frac{1}{a} \sum_{i} \log m_{ij} - \overline{\mu}$$
(2.6)

-4-

$$(\overline{\alpha\beta})_{ij} = \log m_{ij} - \frac{1}{a} \sum_{i} \log m_{ij} - \frac{1}{b} \sum_{j} \log m_{ij} + \overline{\mu}.$$

We note that models (2.4) and (2.5) are of the same form, except for the error term. The  $\bar{\alpha}_i$ 's,  $\bar{\beta}_j$ 's and  $(\bar{\alpha\beta})_{ij}$ 's satisfy

$$\sum_{i} \overline{\alpha}_{i} = \sum_{j} \overline{\beta}_{j} = \sum_{i} (\overline{\alpha\beta})_{ij} = \sum_{j} (\overline{\alpha\beta})_{ij} = 0.$$

Let  $\bar{m}_{ij}$  denote the maximum likelihood estimate of  $\bar{m}_{ij}$  (i = 1,2,...,a; j = 1,2,...,b). Under the hypothesis  $H_0: \Pi_{ij} = \Pi_i \Pi_j$  for all i and j (this is called the hypothesis of independence), where  $\Pi_i = \sum_j \Pi_{ij}$  and  $\Pi_j = \sum_i \Pi_{ij}$ , the maximum likelihood estimates of  $\Pi_i$  and  $\Pi_j$  are  $n_i / n_i$  and  $n_i / n_i$ , respectively, where  $n_i = \sum_j n_{ij}$  and  $n_{ij} = \sum_i n_{ij}$ . Hence,  $\hat{m}_{ij} = n_i \cdot n_{ij} / n_i$ . This is the case of proportional subclass frequencies. The corresponding test statistic is

$$x^{2} = \sum_{i,j} (n_{ij} - \hat{m}_{ij})^{2} / \hat{m}_{ij},$$

which under H has an asymptotic chi-squared distribution with  $\theta = (a-1)(b-1)$  degrees of freedom. If  $c^2 = X^2/n_{...}$ , then

$$\phi(D) = \frac{1}{1+c^2}$$
(2.7)

is a measure of departure from proportionality of the subclass frequencies with  $\phi(D)$  attaining the value one when these frequencies are proportional. In the latter case, model (2.5) takes the additive form

$$\log m_{ij} = \bar{u} + \bar{\alpha}_{i} + \bar{\beta}_{j} . \qquad (2.8)$$

Under the hypothesis of complete balance, namely,  $H_0: \pi_{ij} = 1/(ab)$  for all i and j,  $\hat{m}_{ij} = n_{ij}/(ab)$ , and the corresponding statistic,

$$X^{2} = \sum_{i,j}^{\sum} \frac{\left[ \frac{n_{ij} - n_{..} / (ab)}{n_{..} / (ab)} \right]^{2}}{n_{..} / (ab)}, \qquad (2.9)$$

is asymptotically distributed as a chi-squared variate with  $\theta$  = ab - 1 degrees of freedom. A measure of departure from complete balance is then given by

$$\phi(D) = \frac{1}{1+c^2}$$
, (2.10)

where  $c^2 = X^2/n_{1}$ . In this case model (2.8) is reduced to just

 $\log m = \overline{\mu}$ .

# 2.3 The Three-Way Classification Model

Suppose we consider the model

 $y_{ijkl} = \mu + \alpha_{i} + \beta_{j} + \gamma_{k} + (\alpha\beta)_{ij} + (\alpha\gamma)_{ik} + (\beta\gamma)_{jk} + (\alpha\beta\gamma)_{ijk} + \varepsilon_{ijkl} \quad (2.11)$ (i = 1 2,...,a; j = 1,2,...,b; k = 1,2,...,c; l = 1,2,...,n\_{ijk}),  $\alpha_{i}$ ,  $\beta_{j}$ , and  $\gamma_{k}$  can be either fixed or random. The design D consists of the cell frequencies,  $n_{111}$ ,  $n_{112}$ , ...,  $n_{abc}$ . Following the approach used in the earlier two models, if  $m_{ijk} = E(n_{ijk})$ , then  $\log m_{ijk}$  can be expressed in terms of the loglinear model  $\log m_{ijk} = \overline{\mu} + \overline{\alpha}_{i} + \overline{\beta}_{j} + \overline{\gamma}_{k} + (\overline{\alpha\beta})_{ij} + (\overline{\alpha\gamma})_{ik} + (\overline{\beta\gamma})_{jk} + (\overline{\alpha\beta\gamma})_{ijk}$ , (2.12) where

$$\sum_{i} \vec{\alpha}_{i} = \sum_{j} \vec{\beta}_{j} = \sum_{k} \vec{\gamma}_{k} = \sum_{i} (\vec{\alpha \beta})_{ij} = \sum_{j} (\vec{\alpha \beta})_{ij} = \dots = \sum_{k} (\vec{\alpha \beta \gamma})_{ijk} = 0$$

From (2.12) several reduced models may be considered. These models are given in Table 1. The goodness-of-fit of these models can be checked by using Pearson's approximate chi-squared statistic

$$x^{2} = \sum_{i,j,k} (n_{ijk} - m_{ijk})^{2} / m_{ijk}, \qquad (2.13)$$

where m is the maximum likelihood estimate of m ijk, or by the likelihood ratio statistic

$$G^{2} = 2 \sum_{i,j,k}^{n} n_{ijk} \log(n_{ijk}/m_{ijk})$$
(2.14)

(see Agresti 1984, p. 48). Both  $X^2$  and  $G^2$  are asymptotically distributed as chi-squared variates with the same degrees of freedom. The  $m_{ijk}$  estimates for the models in Table 1 are given in the same table along with the corresponding

Table 1

Some Loglinear Models For a Three-Factor Experiment

Degrees of Freedom	O	b(a-1)(c-1)	(ab-1)(c-1)	abc∽a∽b-c+2	abc-1
c <sup>2</sup>	$G_1^2 = 0$	62 2	6.3 3	ъ С С	6 2 5
x <sup>2</sup>	$x_1^2 = 0$	x 2 2	x <sup>2</sup> 3	x 4 4	x 5
n tjk**	<sup>n</sup> ij k	n.j. <sup>n</sup> .jk	n.i.k	nk n	abc
Model	I. $\log m_{1jk} = \overline{\mu} + \overline{\alpha}_i + \overline{\beta}_j + \overline{\gamma}_k + (\overline{\alpha}\overline{\beta})_{1j}$ + $(\overline{\alpha}\overline{\gamma})_{1k} + (\overline{\beta}\overline{\gamma})_{jk} + (\overline{\alpha}\overline{\beta}\gamma)_{1jk}$	II <sup>*</sup> log $m_{1\mathbf{j}\mathbf{k}} = \overline{\mu} + \overline{\alpha}_{\mathbf{i}} + \overline{\beta}_{\mathbf{j}} + \overline{\gamma}_{\mathbf{k}} + (\overline{\alpha}\overline{B})_{1\mathbf{j}}$ + $(\frac{\overline{\rho}\gamma}{2})_{\mathbf{j}\mathbf{k}}$	III. log m <sub>ijk</sub> = $\overline{\mu} + \overline{\alpha}_i + \overline{\beta}_j + \overline{\gamma}_k + (\overline{\alpha}\overline{\beta})_{1j}$	IV. log $m_{ijk} = \overline{\mu} + \overline{\alpha}_i + \overline{\beta}_j + \overline{\gamma}_k$	ν. log m <sub>ijk</sub> = μ¯

\* Two other models of similar form exist. \*\* The dots used in the elements of this column denote summation over the corresponding missing subscripts.

-7 -

 $x^2$  and  $G^2$  statistics and associated degrees of freedom. The  $G^2$  statistic has the desirable feature of being monotone increasing as terms are deleted from the full model in (2.12), that is,  $0 = G_1^2 \le G_2^2 \le G_3^2 \le G_4^2 \le G_5^2$  (see Agresti 1984, p. 57). The  $G^2$  statistic can, therefore, be used to compare two nested models (that is, one model is obtained from the other by deleting one or more terms) that give adequate fits to the cell frequencies. Thus, with the help of the  $G^2$  statistic it is possible to identify one or more models in Table 1 that provide adequate fits. For such models departures of cell frequencies from their expected values can be measured by means of the function  $\phi(D)$  in (2.10) where  $c^2$  is given by the corresponding value of  $X^2$  in Table 1 divided by n

Model V in Table 1 corresponds to the case of complete balance, whereas Model IV is associated with the case of proportional subclass frequencies. Model III corresponds to a case of conditional proportional subclass frequencies involving values of i, k for a fixed j, and values of j and k for a fixed i. In Model II we have a case of conditional proportional subclass frequencies involving only values of i and k for a fixed j. Model I is the full loglinear model.

# 2.4 The Two-Fold Nested Classification Model

Let us now consider the model

$$y_{ijk} = \mu + \alpha_i + \beta_{ij} + \epsilon_{ijk}$$
(2.15)

(i = 1,2,...,a; j = 1,2,...,b<sub>i</sub>; k = 1,2,...,n<sub>ij</sub>),  $\alpha_i$  denotes the nesting effect and  $\beta_{ij}$  denotes the nested effect. The design D consists of the values of  $b_1$ ,  $b_2$ , ...,  $b_a$  in addition to the  $n_{ij}$  values. In the complete balance case  $b_i = b_i$ for all i and  $n_{ij} = n$  for all i and j. A condition weaker than complete balance is last-stage uniformity which requires that  $n_{ij} = n$  for all i and j. If. however,  $n_{ij} = n_{ij}$ , for  $j \neq j'$  and i = 1, 2, ..., a, then the design is partially balanced. It is known that when all the effects in (2.15) are random, last-stage

-8-

uniformity is a sufficient condition for the sums of squares, in the conventional analysis of variance table, to be independently distributed as scaled chisquared variates (see Tietjen 1974, p. 575). Under partial balance, however, the sums of squares for the  $\alpha_i$  and  $\beta_{ij}$  effects are independent, but do not have the scaled chi-squared distribution (see Cummings 1972). It is, therefore, of interest to measure departures from complete balance, last-stage uniformity, and partial balance.

The loglinear model corresponding to model (2.15) can be obtained as follows: let  $m_{ij}$  denote the expected frequency  $E(n_{ij})$ . Then,  $m_{ij} = n \cdot \prod_{ij} \prod_{j=1}^{n} \prod_{ij=1}^{n} \prod_{j=1}^{n} \prod_{ij=1}^{n} \prod_{j=1}^{n} \prod_{ij=1}^{n} \prod_{j=1}^{n} \prod_{ij=1}^{n} \prod_{j=1}^{n} \prod_{ij=1}^{n} \prod_{j=1}^{n} \prod_{j=$ 

$$\log m_{ij} = \bar{\mu} + \bar{\alpha}_{i} + \bar{\beta}_{ij} , \qquad (2.16)$$

where

$$\widehat{\mu} = \log n \cdot \cdot + \frac{1}{b} \cdot \frac{a}{i=1} b_i \log \Pi_i + \frac{1}{b} \cdot \frac{a}{\sum} \sum_{i=1}^{b_i} \log \Pi_j | i$$

$$\widehat{\alpha}_i = \log \Pi_i + \frac{1}{b_i} \sum_{j=1}^{b_i} \log \Pi_j | i - \frac{1}{b} \cdot \frac{a}{\sum} b_i \log \Pi_i - \frac{1}{b} \cdot \frac{a}{\sum} \sum_{i=1}^{b_i} \log \Pi_j | i$$

$$\widehat{\beta}_{ij} = \log \Pi_j | i - \frac{1}{b_i} \sum_{i=1}^{b_i} \log \Pi_j | i$$

In the partial balance case,  $\pi_{j|i} = 1/b_i$  for all i and j, hence, the maximum likelihood estimate of  $m_{ij}$  is  $\hat{m}_{ij} = n_{i}./b_i$ , where  $n_{i} = \Sigma_{j=1}^{b_i} n_{ij}$  (i = 1,2,...,a), since in this case  $\hat{\pi}_i = n_{i}./n_{ij}$ . A measure of partial balance is then given by

$$\phi(D) = \frac{1}{1+c^2}$$
(2.17)

-9-

where

$$c^2 = x^2/n_{\star}$$

and

$$x^{2} = \sum_{i,j}^{\infty} \frac{(n_{ij} - n_{i} / b_{j})^{2}}{n_{i} / b_{i}} . \qquad (2.18)$$

Under partial balance  $X^2$  has asymptotically the chi-squared distribution with b.-a degrees of freedom, where  $b_i = \sum_{i=1}^{a} b_i$ . This follows from the fact that in the general case, the number of linearly independent  $\Pi_{ij}$ 's is b.-1 whereas under partial balance this number is just a-1. We note that this case can be represented by the loglinear model

$$\log m_{ij} = \bar{\mu} + \bar{\alpha}_{i} . \qquad (2.19)$$

Under last-stage uniformity,  $\Pi_{ij} = 1/b$  for all i and j. The loglinear model in this case has the form

$$\log m_{ij} = \overline{\mu} . \qquad (2.20)$$

The maximum likelihood estimate of  $m_{ij}$  is given by  $n_{..} = n_{..}/b_{..}$  Hence, a measure of departure from last-stage uniformity is given by (2.17), where

$$c^2 = x^2/n$$

and

$$x^{2} = \sum_{i,j}^{\infty} \frac{(n_{ij} - \bar{n}_{..})^{2}}{\bar{n}_{..}}, \qquad (2.21)$$

which has the asymptotic chi-squared distribution with b -1 degrees of freedom.

Unlike the former two cases, departure from complete balance can be attributed to variation in the values of  $b_1, b_2, \dots, b_a$ , or to variation in the  $n_{ij}$ values. We thus need to measure imbalance with regard to the  $b_i$ 's and also with regard to the  $n_{ij}$ 's. We shall consider that the  $b_i$ 's form a multinomial distribution independently of the multinomial distribution of the  $n_{ij}$ 's, with  $b_i$  being distributed as a binomial  $B(b_i, \tau_i)$ . Hence,  $d_i = E(b_i) = b_i \tau_i$  (i = 1,2, ...,a). A measure of imbalance concerning the b,'s is, therefore, given by

$$\phi_1(D) = \frac{1}{1+c_1^2}$$
, (2.22)

where

 $c_1^2 = X_1^2/b$ .

and

$$x_{1}^{2} = \sum_{i=1}^{a} \frac{(b_{i} - \bar{b}_{i})^{2}}{\bar{b}_{i}}, \qquad (2.23)$$

where  $\vec{b} = b/a$ . This statistic has the asymptotic chi-squared distribution with a-1 degrees of freedom when  $\tau = 1/a$  (i = 1,2,...,a). On the other hand, a measure of imbalance concerning the n 's is

$$\phi_2(D) = \frac{1}{1+c_2^2}$$
, (2.24)

where

 $c_2^2 = x_2^2/n$ ...

and

$$x_{2}^{2} = \sum_{i,j}^{z} \frac{(n_{ij} - \bar{n}_{..})^{2}}{\bar{n}_{..}} . \qquad (2.25)$$

The statistic  $X_2^2$  is the same as the one used in last-stage uniformity. Since the multinomial distribution of the b,'s is independent of the multinomial distribution of the  $n_{ij}$ 's,  $x_1^2$  is statistically independent of  $x_2^2$ , hence  $x_1^2 + x_2^2$ is asymptotically distributed as a chi-squared variate with b +a-2 degrees of freedom. Now, to measure departure from complete balance we use the measure

$$\phi(D) = \frac{1}{1+c^2}$$
, (2.26)

where

$$c^2 = c_1^2 + c_2^2$$
 (2.27)

#### 2.5 The Three-Fold Nested Classification Model

In this section we consider the model

$$y_{ijkl} = \mu + \alpha + \beta_{ij} + \gamma_{ijk} + \varepsilon_{ijkl} , \qquad (2.28)$$

 $(i = 1, 2, ..., a; j = 1, 2, ..., b_i; k = 1, 2, ..., c_{ij}; l = 1, 2, ..., n_{ijk})$ . The values of  $b_i$ ,  $c_{ij}$ , and  $n_{ijk}$  make up the design D. Here different types of balance can be considered; each is a stronger type of balance than the one preceding it:

- i) Last-stage partial balance, that is, partial balance with respect to the n<sub>ijk</sub> values. There are two kinds of such partial balance; in the first kind n<sub>ijk</sub> depends on i and j only, and in the second kind n<sub>ijk</sub> depends on i only.
- ii) Last-stage uniformity when the n ijk's are equal for all values of i,j,
   and k.
- iii) Last-stage uniformity and next-to-last-stage partial balance, that is, when  $c_{ii}$  depends on i only.
  - iv) Last-stage uniformity as well as next-to-last-stage uniformity, that is, when the c<sub>ij</sub>'s are equal for all values of i and j.
  - v) Complete balance. This occurs when equality of frequencies occurs within all the subclasses.

Each type can be characterized by one or more loglinear models and a corresponding measure of imbalance can be obtained accordingly. For example, for Type (iii), if the  $n_{ijk}$ 's are considered to have a multinomial distribution with  $n_{ijk}$  being distributed as a binomial  $B(n_{\cdots},\pi_{ijk})$ , where  $n_{\cdots} = \sum_{i,j,k}^{n} n_{ijk}$ 

-12-

then  $m_{ijk} = E(n_{ijk}) = n_{...}/c_{...}$  under last-stage uniformity. Furthermore, if the  $c_{ij}$ 's have a multinomial distribution, independently of the  $n_{ijk}$ 's, with  $c_{ij}$  distributed as  $B(c_{...,\tau_{ij}})$ , where  $c_{...} = \sum_{i,j} c_{ij}$ , then  $d_{ij} = E(c_{ij}) = c_{...\tau_{i}/b_{i}}$ (i = 1, 2, ..., a), since under partial balance with respect to the  $c_{ij}$ 's,  $\tau_{ij} = \tau_{i}\tau_{j|i} = \tau_{i}/b_{i}$ . Thus, the associated loglinear models for  $m_{ijk}$  and  $d_{ij}$ are

 $\log m_{ijk} = \bar{\mu}_1 ,$  $\log d_{ij} = \bar{\mu}_2 + \bar{\alpha}_i .$ 

A measure of imbalance for Type (iii) balance is, therefore, given by

$$\phi(D) = \frac{1}{1+c^2} ,$$

where

$$c^2 = c_1^2 + c_2^2$$

and where

$$\frac{2}{1} = \frac{1}{n_{\dots}} \sum_{i,j,k} \frac{(n_{ijk} - n_{\dots}/c_{\dots})^2}{n_{\dots}/c_{\dots}}, \qquad (2.29)$$

$$c_{2}^{2} = \frac{1}{c_{\cdot\cdot}} \sum_{i,j}^{c_{ij} - c_{i} / b_{j}^{2}} , \qquad (2.30)$$

since in this case the maximum likelihood estimate of  $d_{ij}$  is  $d_{ij} = c_{\cdot\cdot} \tau_i / b_i$ =  $c_{\cdot\cdot} (c_{i\cdot} / c_{\cdot\cdot}) / b_i = c_{i\cdot} / b_i$ . We note that  $n_{\cdot\cdot} c_1^2$  and  $c_{\cdot\cdot} c_2^2$  are distributed independently as asymptotic chi-squared variates with  $c_{\cdot\cdot} - 1$  and  $b_i - 1$  degrees of freedom, respectively.

# 2.6 A Model With A Mixture Of Cross-Classified And Nested Effects

Consider a model involving three factors, A, B, and C, with A and C crossed and B is nested within A. This model is written as

$$\mathbf{y}_{\mathbf{ijkl}} = \mathbf{\mu} + \alpha_{\mathbf{i}} + \beta_{\mathbf{ij}} + \gamma_{\mathbf{k}} + (\alpha_{\gamma})_{\mathbf{ik}} + (\beta_{\gamma})_{\mathbf{ijk}} + \epsilon_{\mathbf{ijkl}}$$
(2.31)

 $(i = 1, 2, ..., a; j = 1, 2, ..., b_i; k = 1, 2, ..., c; l = 1, 2, ..., n_{ijk})$ . Let  $\Pi_{ijk}$  denote the probability of belonging to level i of A, level j of B nested within i, and level k of C. As before, the  $n_{ijk}$ 's are considered to have a multinomial distribution with  $m_{ijk} = E(n_{ijk}) = n_{ijk}$ .

For this model we can have four types of balance:

- i) Proportional subclass frequencies involving the AB subclasses and the levels of factor C, that is,  $\pi_{ijk} = \pi_{ijk} \pi_{ijk}$ .
- ii) Partial balance with respect to the  $n_{ijk}$  values, that is,  $\pi_{ijk} = \pi_i/(b_ic)$ .
- iii) Last-stage uniformity, that is,  $\Pi_{ijk} = 1/(b_c)$  for all i, j, and k.
- iv) Complete balance, that is,  $\tau_i = 1/a$  and  $\Pi_{ijk} = 1/(b_c)$ , where  $\tau_i$  is the i<sup>th</sup> binomial probability associated with the multinomial distribution of the  $b_i$ 's.

Each of the above four types can be represented by a loglinear model. These models are given in Table 2. Furthermore, for each of these four types a measure of imbalance is obtained by using formula (2.3). The value of  $c^2$  in this formula and the degrees of freedom for the corresponding asymptotic chisquared statistics are also given in Table 2.

## 3. Numerical Examples

i) Cummings and Gaylor (1974) used several designs to illustrate the combined effects of dependence and nonchi-squaredness of the analysis of variance mean squares on the size of Satterthwaite's approximate F-test for variance component testing in a two-fold nested model. We shall consider three of these designs which are described in Table 3 and are also represented graphically in Figure 1.

For each of the three designs we measure departures from partial balance

Table 2

Loglinear Models Associated With Four Types Of Balance For Model (2.31)

Туре	Loglinear Model	2 د	Degrees of Freedom
E)	$\log m_{ijk} = \overline{\mu} + \overline{\alpha}_i + \overline{\beta}_{ij} + \overline{\gamma}_k$	$c_{1}^{2} = \frac{1}{n_{\cdots}} \sum_{1,j,k}^{\Sigma} \frac{(n_{jjk} - n_{jj} \cdot n_{\cdots} \cdot k/n_{\cdots})^{2}}{n_{1j} \cdot n_{\cdots} \cdot k/n_{\cdots}}$	(b ~1) (c-1)
(ii)	$\log m_{ijk} = \ddot{\mu} + \ddot{\alpha}_{i}$	$c_{2}^{2} = \frac{1}{n_{*} \cdots 1} \sum_{1,j,k} \frac{\left[n_{1jk} - n_{1,*} / (b_{1}c)\right]^{2}}{n_{1,*} / (b_{1}c)}$	b, c~a
(ii1)	log m <sub>ijk</sub> = <sup>-</sup>	$c_{3}^{2} = \frac{1}{n1} \sum_{i,j,k} \frac{\left[n_{ijk} - n_{i,j} / (b,c)\right]^{2}}{n_{i,j} + n_{i,j} / (b,c)}$	b, c-1
(iv)	$\log m_{1jk} = \tilde{\mu} \text{ and } \log d_1 = \tilde{\mu}^{\prime},$ where $d_1 = E(b_1)$	$c_4^2 = c_3^2 + \frac{1}{b_1} \sum_{i=1}^{2} \frac{(b_1 - b_i)^2}{b_i^2 / a_1}$	b cta-2

. (

-15-

and last-stage uniformity by using formula (2.17) with X<sup>2</sup> being given by (2.18) for partial balance and by (2.25) for last-stage uniformity. We also measure departure from complete balance by applying formulas (2.26) and (2.27). The results are given in Table 4.

Table 3 Designs For a Two-Fold Nested Model

			í		
		1	2	3	4
	b <sub>i</sub>	1	1	4	4
Design 1	<sup>n</sup> ij	1	4	1,1,1,1	4,4,4,4
	b <sub>i</sub>	2	2	2	2
Design 2	<sup>n</sup> ij	1,5	1,5 .	1,5	1,5
	b i	1	2	2	4
vesign 3	<sup>n</sup> ij	1	1,4	1,8	1,2,3,4

## Table 4

Values of  $\phi(D)$  For The Three Designs In Table 3

Design	Partial Balance	Last-Stage Uniformity	Complete Balance
1	1	.735	.58
2	.69	. 69	.69
3	.73	.61	. 54
	-		

From Table 4 we note that other than partial balance for Design 1, none of the designs has strong balance properties. Of all three designs, Design 3 is the most unbalanced with respect to last-stage uniformity and complete balance.



-17-

ii) Bliss (1967, p. 355) described a nested experiment involving three factors with an associated model of the form given by (2.28). In this experiment, a=11 and the design D consists of the following elements:  $b_i=3$  (i = 1,2,...,11);  $c_{i1}=2$ ,  $c_{i2}=2$ ,  $c_{i3}=1$  (i = 1,2,...,11);  $n_{i11}=2$ ,  $n_{i12}=2$ ,  $n_{i21}=1$ ,  $n_{i22}=1$ ,  $n_{i31}=1$  (i = 1,2,...,11). Graphically, for each value of i, the design D can be depicted as in Figure 2.

The design D is partially balanced of the first kind (see Section 2.5). A measure of departure from Type (ii) balance is given by

$$\phi(D) = \frac{1}{1+c_1^2},$$

where  $c_1^2$  is described in (2.29), hence  $\phi(D) = .89$ . The measure for Type (iii) balance is given by

$$\phi(D) = \frac{1}{1+c_1^2+c_2^2}$$
,

where  $c_2^2$  is described in (2.30), hence  $\phi(D) = .83$ . As for Type (iv) balance, the corresponding measure is

$$\phi(D) = \frac{1}{1+c_1^2+c_3^2}$$
,

where

$$c_{3}^{2} = \frac{1}{c_{..}} \sum_{i,j=\bar{c}..}^{(c_{ij}-\bar{c}_{..})^{2}}$$

where  $\bar{c}_{...} = c_{...}/b_{...}$ , hence,  $\phi(D) = .83$ . We note that this is equal to the previous measure value for Type (iii) since both  $c_{i...}$  and  $b_{i...}$  in formula (2.30) do not depend on i, thus,  $c_{i..}/b_{i...} = c_{...}/b_{...} = \bar{c}_{...}$ . We also note that since the  $b_{i...}$ 's are equal, the value  $\phi(D) = .83$  is also a measure of departure from complete balance, which is Type (v) balance.



#### 4. Concluding Remarks

We have introduced a procedure for measuring the degree of imbalance that is associated with an unbalanced model. The procedure applies to cross classification models, nested classification models, and to models with a mixture of cross-classified and nested effects. It can also be used to measure departures from different types of balance, especially in nested models where imbalance can affect various stages of the nested design. Several examples of unbalanced models were studied. From these examples it is easy to see that this procedure is general enough to apply to any unbalanced model.

With the help of this procedure it is now possible to describe in a quantitative manner different kinds of imbalance, such as extreme imbalance, moderate imbalance, and near balance. This can serve as an indicator of the suitability of the approximate methods that are adapted from balanced-data-based procedures and used to analyze an unbalanced model, particularly, when the appropriate measure value is near unity. It is to be cautioned, however, that low values of that measure do not necessarily mean that such approximate methods are inadequate. Cummings and Gaylor (1974), for example, noted that for some extremely unbalanced design, namely, Design 3 in Table 3, their approximate F-test performed very well. They attributed this behavior to counterbalancing effects which appear to reduce, rather than compound, the effect of imbalance on the standard analysis of variance.

## References

Agresti, A., 1984: <u>Analysis of Ordinal Categorical Data</u>. Wiley, New York. Ahrens, H. and R. Pincus, 1981: On two measures of unbalancedness in a one-way model and their relation to efficiency. <u>Biom.</u> J. 23, 227-235.

Bliss, C.I., 1967: Statistics in Biology, Volume 1. McGraw-Hill, New York.

-20-

Cummings, W.B., 1972: Variance component testing in unbalanced nected designs. North Carolina State University Institute of Statistics Mimeo Series No. 843.

Cummings, W.B. and D.W. Gaylor, 1974: Variance component testing in unbalanced nested designs. J. Amer. Statist. Assoc. 69, 765-771.

Hess, J.L., 1979: Sensitivity of MINQUE with respect to a priori weights. Biometrics 35, 645-649.

Tietjen, G.L., 1974: Exact and approximate tests for unbalanced random effects designs. Biometrics 30, 573-581.

Author's address:

E. S. A

Prof. Andre' I. Khuri Department of Statistics Nuclear Sciences Center The University of Florida Gainesville, Florida 32611 U.S.A.

