# Log-linear Modeling Under Generalized Inverse Sampling Scheme

**Soumi Lahiri**[1] **and Sunil Dhar**[2]

[1] Department of Mathematical Sciences
New Jersey Institute of Technology
University Heights, Newark, NJ 07102

[2] Department of Mathematical Sciences
New Jersey Institute of Technology
University Heights, Newark, NJ 07102

# Log-linear Modeling Under Generalized Inverse Sampling Scheme

SOUMI LAHIRI AND SUNIL K. DHAR

Department of Mathematical Sciences, New Jersey Institute of
Technology, Newark, NJ-07111, USA

## SUMMARY

This paper discusses the log-linear model for multi-way contingency table, where the cell values represent the frequency counts that follow an extended negative multinomial distribution. This is an extension of negative multinomial log-linear model described by Evans (1989). The parameters of the new model are estimated by maximum likelihood method. The likelihood ratio test for the general log-linear hypothesis is also derived. A practical application of the log-linear model under the generalized inverse sampling scheme has also been demonstrated by an example.

## 1. INTRODUCTION

Medical and biological researches commonly involve discrete multivariate models. Log-linear models analyze frequency count data. A broad range of sampling plans may arise in biological modeling. Poisson and multinomial samplings are example of direct sampling methods. These sampling models assume, independent cell counts and negatively correlated cell counts, respectively. Moreover, Poisson regression models can only be used where sample mean and sample variance are almost equal. However, in reality, quite often the sample variance is either larger than the sample mean, a case of over dispersion, or the sample variance is smaller than the sample mean, a case of under dispersion. Also, the cell counts for some models can be positively correlated, or sometimes direct sampling methods are not realistic for scientific reasons. In these cases there is a need for inverse sampling methods, e.g., the negative multinomial model. Inverse sampling method is a sampling plan where observations are taken from a population until a predetermined number of "success" is obtained. It is usually used to draw

inference about a rare event. Extended negative multinomial sampling is a generalized inverse sampling scheme, Dhar (1995). It is used when the population consists of more than one rare event and a predetermined number of the rare events are observed.

The test procedures used for direct sampling schemes such as Poisson or Multinomial sampling are not valid under inverse sampling schemes such as negative multinomial sampling (Bishop et al., 1975, p. 455). Therefore, for extended negative multinomial sampling, Steyn (1955, 1959) first gave a Pearson type chi-square test for independence in $R \times C$ contingency tables, where the cell frequency followed a negative multinomial distribution. Bonett (1985a, 1985b) applied the method of minimum chi-square to obtain parameter estimates in negative multinomial log-linear and logit models. He also deduced a Wald test for the general log-linear hypothesis under inverse sampling scheme. Evans and Bonett (1989) presents the maximum likelihood estimator of the negative multinomial log-linear model parameters, giving closed form of the likelihood ratio test statistic for the linear constraints of the regression parameters.

The maximum likelihood estimation method and the likelihood ratio test for the extended negative multinomial log-linear model are presented extending the work of Bonett and Evan (1989) and their earlier results. In Section 2, we define the log-linear model under generalized inverse sampling scheme. Maximum likelihood estimator of the model parameters are derived in Section 3. Section 4 gives the test statistic for the general log-linear model and Section 5 describes the application of this new model.

## 2. EXTENDED NEGATIVE MULTINOMIAL LOG-LINEAR MODEL

Consider a sequence of independent trials as in Dhar (1995), where one of the events $A_i$ occurs with probability $p_i$, $i = -r, \cdots, -1, 1, \cdots, n$, $\sum_{i=-r, i\neq 0}^{n} p_i = 1$. Suppose that $A_{-r}, A_{-(r-1)}, \cdots, A_{-1}$ are the rare events. Let $f_i$ represent the frequency with which $A_i$ occurs until we get a total of $k$ (predetermined value) observations of at least one of the $A_i$'s , $i \in -r, \cdots, -1$. Then the distribution of $\mathbf{f} = (f_{-r}, \cdots, f_{-1}, f_1, \cdots, f_n)'$ is said to follow an extended negative multinomial distribution with parameters $k$ and $\mathbf{p} = (p_{-r}, \cdots, p_{-1}, p_1, \cdots, p_n)'$ with the joint probability density

2

function given as

$$\frac{(\sum_{i=1}^{n} f_i + k - 1)!}{\prod_{i=1}^{n} f_i! \, (k-1)!} \; (\textstyle\sum_{i=1}^{r} p_{-i})^k \; p_1^{f_1} ... \, p_n^{f_n} \; \frac{k!}{\prod_{i=1}^{r} f_{-i}!} \; (p_1^*)^{f_{-1}} ... \, (p_r^*)^{f_{-r}}, \quad (1)$$

where $p_i^* = \dfrac{p_{-i}}{\sum_{i=1}^{n} p_{-i}}, \quad i = 1, ...r, \quad \sum_{i=-r, i \neq 0}^{n} p_i = 1$ and $k = \sum_{i=1}^{r} f_{-i}$. Here, $'$
denotes transpose of a matrix.

The mean vector $\boldsymbol{\mu}$ of $\mathbf{f}$ is a $(n+r) \times 1$ vector and the dispersion matrix
of $\mathbf{f}$ is a $(n+r) \times (n+r)$ block diagonal matrix $\boldsymbol{\Sigma}_f$ of rank $(n+r)$. Both are
computed by using moment generating function (m.g.f.) method and have
the following form:

$$\boldsymbol{\mu} \;\; = \;\; (\mu_{-r}, \cdots, \mu_{-1}, \mu_1, \cdots, \mu_n) = k(\sum_{i=1}^{r} p_{-i})^{-1} \mathbf{p};$$

$$\boldsymbol{\Sigma}_f = \begin{pmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_2 \end{pmatrix}, \quad \text{where}$$

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} kp_1^*(1-p_1^*) & -kp_1^* p_2^* & \cdots & -kp_1^* p_r^* \\ -kp_1^* p_2^* & kp_2^*(1-p_2^*) & \cdots & -kp_2^* p_r^* \\ \vdots & \vdots & \ddots & \vdots \\ -kp_1^* p_r^* & -kp_2^* p_r^* & \cdots & kp_r^*(1-p_r^*) \end{pmatrix},$$

with $p_i^*$ as in equation (1)

and

$\boldsymbol{\Sigma}_2 = ((\boldsymbol{\mu_1}\boldsymbol{\mu_1}')/k + \mathbf{D}_{\boldsymbol{\mu_1}}),$

with $\boldsymbol{\mu_1} = (\mu_1, \mu_2, \cdots, \mu_n)$ and $\mathbf{D}_{\boldsymbol{\mu_1}}$ as the diagonal matrix with elements
of $\boldsymbol{\mu_1}$ along the diagonals.

3

The extended negative multinomial log-linear model is defined as

$$\mathbf{f} = \exp(\mathbf{X}\boldsymbol{\beta}) + \boldsymbol{\delta}, \tag{2}$$

where $\mathbf{X}$ is a $(n+r) \times q$ $(q \leq n+r)$ full rank design matrix, consisting of intercept, main effects and interaction effects, $\boldsymbol{\beta}$ is a $q \times 1$ vector of unknown parameters, and $\boldsymbol{\delta}$ is a $t \times 1$ random error vector with $\mathrm{E}(\boldsymbol{\delta})=0$. The notation 'exp' of a vector in (2) means exponential applied to each component. Here, $\sum_{i=1}^{r} f_{-i}$ is assumed to be a predetermined constant, say k, and $\mathbf{f}$ follows an extended negative multinomial distribution with parameters k and $(n+r) \times 1$ vector $\boldsymbol{\mu} = \mathrm{E}(\mathbf{f}) = \exp(\mathbf{X}\boldsymbol{\beta})$.

## 3. MAXIMUM LIKELIHOOD ESTIMATION OF THE MODEL PARAMETERS

The likelihood function of the extended negative multinomial distribution can be written in the following closed form. In order to express this closed form the following notations are used, $N^* = \mathbf{1}'\mathbf{f}$ and $N = \mathbf{1}'\boldsymbol{\mu}$, where $\mathbf{1}$ is the vector of $(n+r)$ ones. The kernel of the extended negative multinomial log-likelihood function is given by

$$
\begin{aligned}
L(\boldsymbol{\beta}) &= \sum_{i=-r}^{n} f_i \, ln(\mu_i) - (k + \sum_{i=1}^{n} f_i) \, ln(\frac{k}{\sum_{i=1}^{r} p_{-i}}) \\
&= \mathbf{f}' \, ln(\boldsymbol{\mu}) - (\mathbf{1}' \, \mathbf{f}) \, ln(\frac{kp_{-r}}{\sum_{i=1}^{r} p_{-i}} + \cdots + \frac{kp_{-1}}{\sum_{i=1}^{r} p_{-i}} \\
&\quad + \frac{kp_1}{\sum_{i=1}^{r} p_{-i}} + \cdots + \frac{kp_n}{\sum_{i=1}^{r} p_{-i}}), \text{where} \sum_{i=-r, i\neq 0}^{n} p_i = 1 \text{ is used.} \\
&= \mathbf{f}' \, ln(\boldsymbol{\mu}) - (\mathbf{1}' \, \mathbf{f}) \, ln(\mu_{-1} + \cdots + \mu_{-r} + \mu_1 + \cdots + \mu_n) \\
&= \mathbf{f}' \, ln(\boldsymbol{\mu}) - (\mathbf{1}' \, \mathbf{f}) \, ln(\mathbf{1}'\boldsymbol{\mu}) \\
&= \mathbf{f}' \, ln(\boldsymbol{\mu}) - N^* \, [ln(N)]. \tag{3}
\end{aligned}
$$

The maximum likelihood estimator (MLE) of $\boldsymbol{\beta}$ can be obtained by maximizing the expression (3) under the constraint $k = \sum_{i=1}^{r} f_{-i}$. The MLE of $\boldsymbol{\beta}$ cannot be expressed in a closed form due to the complex structure of (3), but can be obtained using some iterative method, say, Newton Raphson Algorithm or EM algorithm. Now, the Newton Raphson algorithm requires

the first and second order derivatives of $L(\boldsymbol{\beta})$ with respect to $\boldsymbol{\beta}$. Applying the methods of matrix derivatives from Dwyer (1967) and using term by term partial differentiation with respect to $\beta_i$, the first and second order partial derivatives can be written in the following form:

$$\frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'} = \mathbf{X}'[\mathbf{f} - (N^*/N)\boldsymbol{\mu}] \tag{4}$$

and

$$\frac{\partial^2 L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = -(N^*/N)[\mathbf{X}'(\mathbf{D} - \frac{\boldsymbol{\mu}\boldsymbol{\mu}'}{N})\mathbf{X}], \tag{5}$$

where $\mathbf{D}$ is a diagonal matrix with elements of $\boldsymbol{\mu}$ along diagonals.

The expression of $L(\boldsymbol{\beta})$ and its first and second order partial derivatives are structurally same as those obtained by Evans and Bonett (1989).

One of the popular methods for finding MLE under constraint is penalty function method. The penalty function is defined as

$$A(\boldsymbol{\beta}) = c(\mathbf{k}_1'\boldsymbol{\mu} - \mathbf{k}_1'\mathbf{f})^2, \tag{6}$$

where c is an arbitrary large positive constant, known as penalty and $\mathbf{k}_1$ is a $(n+r) \times 1$ vector with 1 in the first r positions and 0 in the remaining positions. So the objective function to maximize is

$$M(\boldsymbol{\beta}) = L(\boldsymbol{\beta}) + A(\boldsymbol{\beta}). \tag{7}$$

The first and second order partial derivatives of $M(\boldsymbol{\beta})$ have the following closed form

$$M'(\boldsymbol{\beta}) = \frac{\partial M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'} = \mathbf{X}'[(\mathbf{f} - (N^*/N)\boldsymbol{\mu}) + 2c\, \mathbf{D}\mathbf{k}_1(\mathbf{k}_1'\boldsymbol{\mu} - k)]$$

$$M''(\boldsymbol{\beta}) = \frac{\partial^2 M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = -\mathbf{X}'[(N^*/N)(\mathbf{D} - \frac{\boldsymbol{\mu}\boldsymbol{\mu}'}{N}) - 2c\, \mathbf{D}\mathbf{k}_1\mathbf{k}_1'(2\mathbf{D} - k)]\mathbf{X}.$$

Using the Newton Raphson algorithm, the MLE of $\boldsymbol{\beta}$ can be obtained iteratively as

$$\begin{aligned} \mathbf{b}_{m+1} &= \mathbf{b}_m - [M''(\mathbf{b}_m)]^{-1} M'(\mathbf{b}_m) \\ &= \mathbf{b}_m + \mathbf{P}_m \mathbf{g}_m, \end{aligned} \tag{8}$$

where $\mathbf{P}_m = [\mathbf{X}'(N^*/N)(\mathbf{D}_m - \frac{\boldsymbol{\mu}_m \boldsymbol{\mu}_m'}{N}) - 2c\, \mathbf{D}_m \mathbf{k}_1 \mathbf{k}_1'(2\mathbf{D}_m - k)\mathbf{X}]^{-1}$,
$\mathbf{g}_m = \mathbf{X}'[(\mathbf{f} - (N^*/N)\boldsymbol{\mu}_m) + 2c\, \mathbf{D}_m \mathbf{k}_1(\mathbf{k}_1'\boldsymbol{\mu}_m - k)]$, $\boldsymbol{\mu}_m = \exp(\mathbf{X}\mathbf{b}_m)$. The diagonal matrix $\mathbf{D}_m$ has the elements of $\boldsymbol{\mu}_m$ along the principal diagonal. The

MLE of $\boldsymbol{\beta}$, denoted as $\widehat{\boldsymbol{\beta}} = \mathbf{b}_{m+1}$, is obtained when the difference between $\mathbf{b}_{m+1}$ and $\mathbf{b}_m$ is arbitrarily small. The initial value $\mathbf{b}_0$ is taken as the least square estimate, $\mathbf{b}_0 = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'ln(\mathbf{f})$, setting $ln(0) = 0$. The invariance property of MLE, yields MLE of $\boldsymbol{\mu}$ to be $\widehat{\boldsymbol{\mu}} = \exp(\mathbf{X}\widehat{\boldsymbol{\beta}})$.

The asymptotic covariance matrix of $\widehat{\boldsymbol{\beta}}$ can be obtained by expanding the expression of (7) by the mean value theorem (MVT) around the true parameter $\boldsymbol{\beta_0}$ as follows:

$$\frac{\partial M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'} = \frac{\partial M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'}|_{(\boldsymbol{\beta}=\boldsymbol{\beta_0})} + (\boldsymbol{\beta} - \boldsymbol{\beta_0})\frac{\partial^2 M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}\partial \boldsymbol{\beta}'}|_{(\boldsymbol{\beta}=\boldsymbol{\beta_1})},$$

where $\boldsymbol{\beta_1}$ lies in the small neighborhood of $\boldsymbol{\beta_0}$. Now letting $\dfrac{\partial M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'} = 0$, gives

$$-\mathbf{X}'[(\mathbf{f} - (N^*/N)\boldsymbol{\mu}) + 2c\,\mathbf{Dk}_1(\mathbf{k}_1'\boldsymbol{\mu} - k)] = \frac{\partial^2 M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}\partial \boldsymbol{\beta}'}|_{(\boldsymbol{\beta}=\boldsymbol{\beta_1})}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta_0})$$

Note that for large sample size $\boldsymbol{\beta_1} \approx \boldsymbol{\beta_0}$ and from the above expression the asymptotic covariance matrix of $\widehat{\boldsymbol{\beta}}$ can be computed as

$$\boldsymbol{\Sigma}_{\widehat{\boldsymbol{\beta}}} = [\frac{\partial^2 M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}\partial \boldsymbol{\beta}'}]^{-1}(\mathbf{X}'\,\boldsymbol{\Sigma}_f\,\mathbf{X})[\frac{\partial^2 M(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}\partial \boldsymbol{\beta}'}]^{-1}.$$

Therefore the estimate of the asymptotic covariance matrix of $\widehat{\boldsymbol{\beta}}$ is given by

$$\widehat{\boldsymbol{\Sigma}}_{\widehat{\boldsymbol{\beta}}} = \mathbf{P}[\mathbf{X}'\,\widehat{\boldsymbol{\Sigma}}_f\,\mathbf{X}]\mathbf{P}', \tag{9}$$

where $\mathbf{P}$ and $\mathbf{D}$ are the values of $\mathbf{P}_m$ and $\mathbf{D}_m$ obtained in the last iteration of (8).

## 4. HYPOTHESIS TESTING

In this section the likelihood ratio test statistic of the general linear hypothesis $H_0 : \mathbf{H}\boldsymbol{\beta} = 0$ versus its negation is derived, where $\mathbf{H}$ is a $p \times q$ known matrix of rank $p$. Evans (1989) derived the likelihood ratio test statistic for the general log-linear hypothesis under negative multinomial sampling. Here the likelihood ratio statistic is computed for the extended multinomial sampling plan. Alternatively, the Wald statistics can also be derived to evaluate the log-linear hypothesis.

Following Graybill (1976, p. 186) and substituting the constraint $\mathbf{H}\boldsymbol{\beta} = 0$ in the model $\ln(\mathbf{f}) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\delta}$, a new reduced model can be obtained as $\ln(\mathbf{f}) = \mathbf{X}(\mathbf{I} - \mathbf{H}^-\mathbf{H})\boldsymbol{\beta} + \boldsymbol{\delta}$, where $\mathbf{H}^-$ denotes the generalized inverse of the matrix $\mathbf{H}$. The likelihood ratio test statistic $\lambda$ is obtained as

$$-2\ln(\lambda) = 2\mathbf{f}'\,\mathbf{X}\,\mathbf{H}^-\mathbf{H}\,\widehat{\boldsymbol{\beta}} + 2N^*\ln\left(\frac{\mathbf{1}'\exp(\mathbf{X}(\mathbf{I} - \mathbf{H}^-\mathbf{H})\widehat{\boldsymbol{\beta}})}{\mathbf{1}'\exp(\mathbf{X}\widehat{\boldsymbol{\beta}})}\right) \quad (10)$$

which asymptotically follows a chi-square distribution with $n + r - q$ degrees of freedom (d.f.).

## 5. EXAMPLE

Oxybutynin is the most commonly used drug for the treatment of overactive bladder symptom. But this drug has several adverse effects, for example, dry mouth, dyspepsia, dysuria, upper respiratory tract infection, lower respiratory tract infection, urinary infection etc. Some of them are so serious that patients even cannot continue the treatment. An alternative of this drug is tolterodine. Our objective is to find out whether tolterodine also has significant serious adverse effects.

Suppose a group of patients reported with overactive bladder problems was given oxybutin and another group was prescribed tolterodine and were asked to report after certain time. Then three variables each with two levels were recorded for each patient: Gender(male or female), used tolterodine (yes or no), and suffering from serious adverse effects (yes or no). Samples were recorded until 15 patients who were prescribed tolterodine reported serious adverse effects. Hypothetical data for this study is given below.

|  |  | Tolterodine Used | | | | |
|---|---|---|---|---|---|---|
|  |  | Yes | | | No | |
|  |  | Serious Adverse Effects | | | Serious Adverse Effects | |
|  |  | Yes | No | | Yes | No |
|  | male | 8 | 30 | | 19 | 12 |
| Gender |  |  |  | |  |  |
|  | female | 7 | 38 | | 25 | 15 |

**Objective:** To find the relationship between the observed counts and three variables (gender, drug used and adverse effects) along with their interactions.

The log-linear model for this example will be

$$\ln \mathbf{f} \;=\; \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\delta},$$

where $\mathbf{X}$ contains three main effects (Tolterodine used, suffering from adverse effects and gender respectively) along with their all possible. Therefore the form of the design matrix $\mathbf{X}$ will be

$$\mathbf{X} \;=\; \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \text{ and } \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \end{pmatrix},$$

where $\beta_0$ = general mean effect, $\beta_1$ = differential effect due to tolterodine, $\beta_2$ = differential effect due to adverse effect, $\beta_3$ = differential effect due to gender, $\beta_4$ = differential effect due to interaction of tolterodine and gender, $\beta_5$ = differential effect due to interaction of tolterodine and adverse effect, $\beta_6$ = differential effect due to interaction of gender and adverse effect. Here $\mathbf{f}$ denotes the frequency counts and follows an extended negative multinomial distribution with parameters $(k = 15, p_{-1}, p_{-2}, p_1, \cdots, p_6)$. Then the maximum likelihood estimates of the model parameters are

$\widehat{\beta}_0 = 2.75, \ \widehat{\beta}_1 = -0.33, \ \widehat{\beta}_2 = 0.43, \ \widehat{\beta}_3 = 0.86, \ \widehat{\beta}_4 = 0.13, \ \widehat{\beta}_5 = -2.00,$ and $\widehat{\beta}_6 = 0.14.$

The following table shows the estimated value of the expected frequency of **f**.

TABLE 1
Estimation of the frequency counts

| Cell | **f** | $\widehat{\boldsymbol{\mu}}$ |
|------|----|---------|
| 111 | 8 | 7.2598 |
| 112 | 19 | 19.7369 |
| 121 | 30 | 30.7356 |
| 122 | 12 | 11.2594 |
| 211 | 7 | 7.7383 |
| 212 | 25 | 24.2579 |
| 221 | 38 | 37.2563 |
| 222 | 15 | 15.7374 |

The sign of $\widehat{\beta}_1$ and $\widehat{\beta}_5$ implies that the use of drug tolterodine and adverse effects due to its use are negatively correlated.

Our objective is to test the null hypothesis that the following two way interactions, the gender by adverse effects, and the gender by tolterodine, are all equal to zero in the above model, that is to test $\mathbf{H}\boldsymbol{\beta} = 0$, where $\mathbf{H} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$. So the reduced model contains only the intercept, the three main effects and the tolterodine by adverse effects interaction. The likelihood ratio statistic follows a chi-squared distribution with 1 d.f. and the value of the statistic equals 0.1881 which suggests that the reduced model is appropriate at 1% level of significance.

# References

[1] Bishop, Y.M.M., Fienberg, S.E., and Holland, P.W. (1975). "Discrete Multivariate Analysis," *Cambridge: MIT Press.*

[2] Bonett, D.G. (1985a). "A linear negative multinomial model," *Statist.& Probability Letters*, **3**, 127-129.

[3] Bonett, D.G. (1985b)."The negative multinomial logit model," *Commun. Statist.-Theory Methods*, **14(7)**, 1713-1717.

[4] Dhar, S.K. (1995)."Extension of a Negative Multinomial Model," *Commun. Statist.-Theory Methods*, **24(1)**, 39-57.

[5] Dwyer, P.S. (1967)."Some applications of matrix derivatives in Multivariate analysis," *J. American Statist. Assoc*, **62**, 607-625.

[6] Evans, M.A. and Bonett, D.G. (1989)."Maximum likelihood estimation for the negative multinomial log-linear model," *Commun. Statist.-Theory Methods*, **18(11)**, 4059-4067.

[7] Graybill, F.A. (1976)."Theory and Application of the Linear Model," *Wadsworth Publishing Company, Inc., Belmont, California 94002.*

[8] Steyn, H.S. (1959)."On $\chi^2$ tests for contingency tables of negative multinomial types," *Statistica Neerlandica*, **13**, 433-444.