Correlation Coefficient between Two types of Children in a Family

Talawar A. S.¹, Sarvade R. M.²

^{1, 2}Department of Statistics, Karnataka University, Dharwad, Karnataka State, India

Abstract: Results show that both the modified negative binomial and negative binomial distribution give good fit. Comparing with Chisquare values, the family size fits better with a modified negative binomial distribution than with a negative binomial distribution. The observed correlation co-efficient between the two types of children is 0.2403 and the estimated correlation co-efficient between boys and girls, in case of negative binomial distribution is 0.216, when $p\neq q$ and 0.2164, when p=q=1/2 and in case of modified negative binomial distribution is 0.2244, when $p\neq q$ and 0.2249, when p=q=1/2.

Keywords: probability model, truncated model, family size, correlation co-efficient, Chi-square test

1. Introduction

Statisticians, demographers and other scientists have been concerned about the family size and the distribution of boys and girls within families of different sizes and also the correlation between the two types of children. Rao et al. (1973) considered family size as a random variable following negative binomial distribution. They developed the more general expression for the correlation between two types of children in a family. Hamdan (1975) obtained the correlation between the two types of children in a family using truncated negative binomial distribution excluding the childless families. Gupta (1976) obtained the correlation coefficient for the modified power series distribution (MPSD). Rao (1979) obtained the correlation between the numbers of two types of children in a family using modified power series distribution. Janardan (1982) obtained the correlation between the numbers of two types of children in a family using the Markov-polya model. In the present paper we have correlation co-efficient for some of the probability models. The purpose of this paper is to give a good fit to the distribution of births considering negative binomial distribution, modified negative binomial and truncated negative binomial and check their adequacy to the data and also to obtain the correlation coefficient between the two types of children in a family using the above distribution. For this purpose we have used cross-sectional data collected from randomly chosen 470 families (with completed fertility), out of which 32 are childless families.

2. Distribution of a Family Size Considering the Families with Childless Families

Let the family size N, be a random variable formed from the two types of children boys (B) and girls (G), that is N=B+G.

(a) Negative Binomial Distribution

Let the family size N has a negative binomial distribution with the probability function

$$P(N=k) = {\binom{r+k-1}{k}} \theta^r (1-\theta)^k, k = 0, 1, 2, \dots (2.1)$$

Where r>0 and 0< θ <1. The mean and variance of the family size, N are given by

$$E(N) = \frac{r(1 - \theta)}{\theta}$$
 and $V(N) = \frac{r(1 - \theta)}{\theta^2}$

The parameters are estimated using the method of maximum likelihood estimation.

(b) Modified Negative Binomial Distribution

Let the family size N has a modified negative binomial distribution with probability function

$$P(N=0) = P_0 = 1 - \alpha + \alpha \theta^r$$

$$P(N=k) = \alpha \binom{r+k-1}{k} \theta^r (1-\theta)^k, k = 1, 2, \dots$$
(2.2)

Where $0 < \theta < 1$, $0 < \alpha < 1$ and α is proportion of fecundability (Sharma, 1995). The mean and variance of the family size N are given by

$$E(N) = \frac{\alpha r(1-\theta)}{\theta} \text{ and } V(N) = \frac{\alpha r(1-\theta)}{\theta^2} \left[1 + (1-\alpha)r(1-\theta)\right]$$

The parameters are estimated using the method of maximum likelihood estimation.

2.1 Distribution of a family size without childless families

Here the childless families are excluded from data, because the childless families do not contribute to the distribution of births and hence are excluded.

(a) Zero-truncated Negative Binomial Distribution (TNBD)

Assuming that the family size N has a zero-truncated negative binomial distribution with probability function

$$P[N=k] = \frac{(r+k-1)!}{(r-1)!k!} \frac{\theta'(1-\theta)^{k}}{1-\theta^{r}}, k = 1, 2, 3 \dots (3.1)$$

The first two non-central moments of the TNBD are $\alpha r(1 - \theta)$

$$\mu_1 = E(N) = \frac{1}{\theta(1-\theta^r)}$$
$$\mu_2' = E(N^2) = \frac{r(1-\theta)(1+r(1-\theta))}{\theta^2(1-\theta^r)}$$

And therefore the variance of the family size is thus obtained $V(N) = \mu_n^2 - \mu_n^2$

$$= \frac{r(1-\theta)}{\theta^2(1-\theta^r)^2} \left(1 - \theta^r \left(1 + r(1-\theta)\right)\right)$$

The parameters θ and r of (3.1) are estimated using the

Volume 3 Issue 9, September 2014

<u>www.ijsr.net</u>

Licensed Under Creative Commons Attribution CC BY

645

= 1

methodology provided by Brass (1958a and 1958b) and Sampford (1955). We use the methodology provided by Sampford (1955). Therefore We have a simplified equation

$$\varphi(w) = m \exp\left\{-\left(m + \frac{s^2}{m} - 1\right)\phi(w)\right\} + \frac{s^2}{m}w$$

Taking the initial values for

$$w = \frac{1}{s^2/m} \phi'(w) = \frac{s^2}{m} - m\left(m + \frac{s^2}{m} - 1\right) \left[\frac{\phi(w) - w}{w(1 - w)}\right] \exp\left\{-\left(m + \frac{s^2}{m} - 1\right)\phi(w)\right\}$$

Thus from Newton-Rapson method, for the ith iteration

$$w_i = w_{i-1} - \frac{\varphi(w)}{\varphi'(w)}\Big|_{w=w_{i-1}}$$

(b) Logarithmic Series Distribution

Let the random variable N, the family size, has a logarithmic series distribution with probability function

$$\mathbb{P}(N=k) = \frac{\alpha \theta^{\kappa}}{\theta}, k = 1, 2, 3, \dots, 0 < \theta < 1 (3.2)$$

where $\alpha = -[\log (1-\theta)]^{-1}$

The individual probabilities are the terms in the series expansion of $-\alpha \log(1-\theta)$. The mean and variance of N are given by

$$E(N) = \frac{\alpha \theta}{1-\theta}$$
 and $V(N) = \frac{\alpha \theta (1-\alpha \theta)}{(1-\theta)^2}$

The parameters are obtained using the likelihood equation αθ

$$x = \frac{1}{(1-\theta)}$$

and the table provided by Patil (1962).

(c) Beta-Geometric Distribution:

If we assume that the family size N has a beta-geometric distribution with the following form of probability function

$$\mathbf{P}(\mathbf{N}=\mathbf{k}) = \mu \frac{\prod_{i=1}^{k-1} [1-\mu+(i-1)\theta]}{\prod_{i=1}^{k} [1+(i-1)\theta]}, \mathbf{k} = 1, 2, \dots (3.3)$$

The mean and variance of model (3.3) are given by

$$E[N] = \left[\frac{1-\theta}{\mu-\theta}\right] \text{ and } Var[N] = \frac{\mu(1-\mu)(1-\theta)}{(\mu-\theta)^2(\mu-2\theta)}$$

The parameters μ and θ are estimated by maximum likelihood method (Weinberg and Gladen, 1986).

(b) Modified Negative binomial distribution

Substituting the mean and variance of the modified negative binomial distribution in (4.1), we get

$$\rho_{BG} = \frac{(pq)^{\frac{1}{2}} \left[\frac{\alpha r(1-\theta)}{\theta^{2}} \left(1 + (1-\alpha)r(1-\theta) \right) - \frac{\alpha r(1-\theta)}{\theta} \right]}{\left[\frac{p\alpha r(1-\theta)}{\theta^{2}} \left(1 + (1-\alpha)r(1-\theta) \right) + \frac{q\alpha r(1-\theta)}{\theta} \right]^{\frac{1}{2}} \left[\frac{q\alpha r(1-\theta)}{\theta^{2}} \left(1 + (1-\alpha)r(1-\theta) \right) + \frac{p\alpha r(1-\theta)}{\theta} \right]^{\frac{1}{2}}}{\rho_{BG}} = \frac{(pq)^{\frac{1}{2}} \left[\frac{\alpha r(1-\theta)}{\theta^{2}} \left(1 + (1-\alpha)r(1-\theta) \right) + \frac{\alpha r(1-\theta)}{\theta} \right]}{\frac{\alpha r(1-\theta)}{\theta^{2}} \left[p(1+(1-\alpha)r(1-\theta)) + q\theta \right]^{\frac{1}{2}} \left[q(1+(1-\alpha)r(1-\theta)) + p\theta \right]^{\frac{1}{2}}}}{\left[p(1+(1-\alpha)r(1-\theta) + q\theta \right]^{\frac{1}{2}} \left[q(1+(1-\alpha)r(1-\theta)) + p\theta \right]^{\frac{1}{2}'}}}$$

Volume 3 Issue 9, September 2014 www.ijsr.net

Licensed Under Creative Commons Attribution CC BY

2. General formula for correlation coefficient:

Let the non-negative integer-valued random variable N denotes for the size of a family and follows any distribution with mean, E(N) and variance, V(N) both finite. Let. We assume that, N=B+G, where B is the number of boys and G is the number of girls in a family for any given size, B is a binomial variable (N, p) where p is the probability of a child being a boy and q=1-p is probability of a child being a girl. The general formula the correlation coefficient between the two types of children B and G due to Rao et al. (1973), is given by

$$\rho_{BG} = \frac{(pq)^{\frac{1}{2}} [V(N) - E(N)]}{[pV(N) + qE(N)]^{\frac{1}{2}} [qV(N) + pE(N)]^{\frac{1}{2}}}, \text{ if } p \neq q$$
$$= \frac{V(N) - E(N)}{V(N) + E(N)}, \text{ if } p = q = \frac{1}{2} (4.1)$$

Where p and q are respectively, the proportions of boys and girls in a family.

(a) Negative binomial distribution

Substituting mean and variance of N for the negative binomial distribution in (4.1) we get

$$\rho_{BG} = \frac{(pq)^{\frac{1}{2}} (1-\theta)}{[(p+q\theta)(q+p\theta)]^{\frac{1}{2}}}, \text{ if } p\neq q$$
$$= \frac{1-\theta}{1+\theta}, \text{ if } p=q=\frac{1}{2} (4.2)$$

Note that the correlation coefficient between B and G (ρ_{BG}) is independent of the parameter r. We get the same ρ_{BG} , when the size of family, N has a geometric distribution.

That is, if N has a geometric distribution with probability function

$$P(N=k) = \theta(1-\theta)^{k}, k=0, 1, 2...$$

Mean and variance of N are
$$E(N) = \frac{1-\theta}{\theta} \text{ and } V(N) = \frac{1-\theta}{\theta^{2}}$$

Therefore, the correlation coefficient between the two types of children, B and G is

8²

$$\rho_{BG} = \frac{(pq)^{\frac{1}{2}} (1-\theta)}{\sqrt{(p+q\theta)(q+p\theta)}}, \text{ if } p \neq q$$
$$= \frac{1-\theta}{1+\theta}, \text{ if } p = q = \frac{1}{2}$$

 $\begin{array}{l} \text{if } p \neq q \\ = \frac{1 + (1 - \alpha)r(1 - \theta) - \theta}{1 + (1 - \alpha)r(1 - \theta) + \theta} \text{ , if } p = q = \frac{1}{2} \end{tabular} \ (4.3) \end{array}$

(c) Truncated Negative binomial distribution:

The correlation coefficient between the two kinds of children, B and G is

$$\rho_{BG} = \frac{(pq)^{\frac{1}{2}} \left[1 - \theta^{r} (1 + r(1 - \theta)) - \theta(1 - \theta^{r})\right]}{\left[p\left(1 - \theta^{r} (1 + r(1 - \theta))\right) + q\theta(1 - \theta^{r})\right]^{\frac{1}{2}} \left[q\left(1 - \theta^{r} (1 + r(1 - \theta))\right) + p\right]^{\frac{1}{2}}}$$

Put $1 \cdot \theta^r = A$ $1 \cdot \theta^r (1 + r(1 \cdot \theta)) \cdot \theta (1 \cdot \theta^r) = 1 \cdot \theta^r \cdot \theta^r r(1 \cdot \theta) \cdot \theta (1 \cdot \theta^r)$ $= A_r \cdot \theta^r (1 \cdot \theta) - \theta A$ $= (1 \cdot \theta) A_r \cdot \theta^r (1 \cdot \theta)$ $= (1 \cdot \theta) (A_r \cdot \theta^r)$ we get $(rec)^{\frac{1}{2}} (1 \cdot \theta) (A_r \cdot \theta^r)$

$$P_{BG} = \frac{(\mathbf{p} + \mathbf{q})}{[(\mathbf{p} + \mathbf{q})\mathbf{A} - \mathbf{pr}(1 - \mathbf{\theta})\mathbf{\theta}^{T}]^{\frac{1}{2}}[(\mathbf{q} + \mathbf{p})\mathbf{A} - \mathbf{qr}(1 - \mathbf{\theta})\mathbf{\theta}^{T}]^{\frac{1}{2}}}, \mathbf{p} \neq \mathbf{q}$$

Which is similar to Hamdan (1975) (equation (2.2)).

(e) Beta-Geometric Distribution:

The correlation coefficient between the two types of children, B and G is

$$\begin{split} \rho_{BG} &= \frac{(pq)^{\frac{1}{2}} \left[\frac{\mu(1-\mu)(1-\theta)}{(\mu-\theta)^2(\mu-2\theta)} - \frac{(1-\theta)}{(\mu-\theta)} \right]}{\left[\frac{p\mu(1-\mu)(1-\theta)}{(\mu-\theta)^2(\mu-2\theta)} + \frac{q(1-\theta)}{(\mu-\theta)} \right]^{\frac{1}{2}} \left[\frac{q\mu(1-\mu)(1-\theta)}{(\mu-\theta)^2(\mu-2\theta)} + \frac{p(1-\theta)}{(\mu-\theta)} \right]^{\frac{1}{2}}}{\left[pq \right]^{\frac{1}{2}} \left[\mu(1-\mu) - (\mu-\theta)(\mu-2\theta) \right]} \\ \rho_{BG} &= \frac{(pq)^{\frac{1}{2}} \left[\mu(1-\mu) - (\mu-\theta)(\mu-2\theta) \right]}{\left[p\mu(1-\mu) + q(\mu-\theta)(\mu-2\theta) \right]^{\frac{1}{2}} \left[q\mu(1-\mu) + p(\mu-\theta)(\mu-2\theta) \right]^{\frac{1}{2}}} \\ &= \frac{\mu(1-\mu) - (\mu-\theta)(\mu-2\theta)}{\mu(1-\mu) + (\mu-\theta)(\mu-2\theta)}, \text{if } p = q = \frac{1}{2} (4.6) \end{split}$$

3. Application of the Models

The following **Table 1** gives the observed numbers of males and females offspring among the 470 families. Total number of boys is 899 and total number of girls is 785. Therefore the total number of children in 470 families is 1684. The proportion of males (p) is 0.534 and the proportion of females (q) is 0.466. Figures in parenthesis give the percentage values of particular combinations of males and females.

Fitting of Probability Models with zeroes and without zeroes

Empirical statistics shows that both the modified negative binomial and negative binomial distribution give good fit. Comparing with Chi-square values, the human family size fits better with a modified negative binomial distribution than with a negative binomial distribution (**Table 2**). The observed correlation co-efficient between the two types of children is 0.2403 and the estimated correlation co-efficient between B and G, in case of negative binomial distribution

$\rho_{BG} = \frac{(1-\theta)(\mathbf{A}-\mathbf{r}\theta^{t})}{[(1+\theta)\mathbf{A}-\mathbf{r}(1-\theta)\theta^{t}]}, p=\mathbf{q} = \frac{1}{2}$ (4.4)

(d) Logarithmic series distribution:

The correlation coefficient between the two types of children is

$$\rho_{BG} = \frac{(pq)^{\frac{1}{2}} \left[\frac{\alpha \theta(1-\theta)}{(1-\theta)^2} - \frac{\alpha \theta}{(1-\theta)} \right]}{\left[\frac{p\alpha \theta(1-\alpha \theta)}{(1-\theta)^2} + \frac{q\alpha \theta}{(1-\theta)} \right]^{\frac{1}{2}} \left[\frac{q\alpha \theta(1-\alpha \theta)}{(1-\theta)^2} + \frac{p\alpha \theta}{(1-\theta)} \right]^{\frac{1}{2}}}$$
$$= \frac{(pq)^{\frac{1}{2}} \theta(1-\alpha)}{\left[p - p\alpha \theta + q - q\theta \right]^{\frac{1}{2}} \left[q - q\alpha \theta + p - p\theta \right]^{\frac{1}{2}}}}{p_{BG}}$$
$$\rho_{BG} = \frac{\frac{(pq)^{\frac{1}{2}} \theta(1-\alpha)}{(1-\theta)(p\alpha + q)\left[\frac{1}{2}(1-\theta)(q\alpha + p)\right]^{\frac{1}{2}}}}{\frac{1}{(1-\theta)(p\alpha + q)\left[\frac{1}{2}(1-\theta)(q\alpha + p)\right]^{\frac{1}{2}}}}, p \neq q, \text{ since } p+q=1$$
$$= \frac{\theta(1-\alpha)}{2-\theta(1+\alpha)}, p = q = \frac{1}{2} (4.5)$$

is 0.216, when $p\neq q$ and 0.2164, when p=q=1/2 and in case of modified negative binomial distribution is 0.2244, when $p\neq q$ and 0.2249, when p=q=1/2.

 Table 1: Observed Numbers of Males and Females offspring among 470 Families

International Journal of Science and Research (IJSR) ISSN (Online): 2319-7064 Impact Factor (2012): 3.358

Number of	Number of female offspring						Total			
male offspring	0	1	2	3	4	5	6	7	8	number of families
9				1						1
8			×		\backslash					0
7			1	A	×	\searrow				2
6	×	7	1	2	A	1				12
5	1	3	4	4	2	2	\setminus	\sum		16
4	4	9	10	5	3	X		\sum	×	32
3	8	18	18	12	×	4	T	1	¥ 1	70
2	25	46	29	19	7	T	1	1	*	129
1	34	37	33	D	11	1	1	2		128
0	32	24	× 15	7	2	\searrow	×	\searrow		80
Total number of families	104	144	111	60	33	10	3	4	1	470

Table 2: Fitting of Probability Models

Size of	Observed	Expected		
family	No. of families	Negative Binomial	Modified NBD	TNBD
		£=6.476 & 0=0.6442	î−6.84 θ̂−0.654 & α̂=0.9858	î − 7.94 & 9 − 0.685
0	32	27	32	-
1	58	63	60	58
2	77	83	81	77
3	94	85	83	94
4	62	71	71	62
5	58	53	53	58
6	33	36	36	33
7	25	23	23	25
8	15	14	14	15
9	8	8	8	8
10	5	4	4	5
11	2	2	2	2
12	1	1	1	1
	470	470	470	438
		$\chi^2 = 4.96^{\rm ns}$	$\chi^2 = 3.86^{\rm ns}$	$\chi^2 = 3.78^{\rm ns}$

^{ns} indicates significant at 5% level of significance

TNBD gives good fit to the size of a family when the childless families are excluded.

References

- [1] Brass, W. (1958a) Simplified methods of fitting the truncated negative binomial distribution, *Biometrika*, 15, 237-250.
- [2] Brass, W. (1958b) The distribution of births in human population, *population studies*, 12, 57-72.
- [3] Gupta R. C. (1976) Application of Modified Power Series Distribution in Genetics, *Sankhya*: The Indian J. of Statistics, Series B, Vol. 38, No2, pp 187-191.
- [4] Hamdan, M. A. (1975) Correlation between the numbers of two types of children when the family size distribution is zero-truncated negative binomial, *Biometrics* 31, 765-766.
- [5] Janardan, K.G., (1982) Correlation between the numbers of two types of children in a family, using the

Volume 3 Issue 9, September 2014

<u>www.ijsr.net</u>

Markov-Polya Model, *Mathematical Biosciences*, 62(1), 123-136.

- [6] Patil, G.P. (1962) Some methods of estimation for the logarithmic series distribution, *Biometrics*, 18, 68-75.
- [7] Rao, B. R., (1981) Correlation between the numbers of two types of children in a family with the mpsd for the family size, *communications in statistics – Theory and methods*, 10(3), 1981.
- [8] Rao, B. R., Mazumdar, S., Walter, H.J., and Li, C.C., (1973) Correlation between the numbers of two types of children in a family, *Biometrics*, 29, 271-279.
- [9] **Sharma, H.L., (1995)** Estimation of parameters involved in the distribution of numbers of two types of children in a family, *Journal of ISPS*, Vol 2, 1-14.

Author Profile

Dr. A. S. Talawar is M.Sc., M.Phil., Ph.D. having experience of 19 years. Presently working as Associate Professor and has published 18 papers in national and international journal

R. M. Sarvade is Ph.D. Student who is about to submit his Thesis shortly. He has three papers published in his name