

## **LINEAR MODELS : ENERGY FORECASTING, PERFORMANCE EVALUATION METHODS.**

**SUGUNA.T**

Research Scholar

M.Phil Mathematics

Bharath Institute Of Higher Education And Research

Mail Id: [rsugu1971@gmail.com](mailto:rsugu1971@gmail.com)

Guide Name: **Dr. R. DEEPA**

Assistant Professor, Department Of Mathematics

Bharath Institute Of Higher Education And Research

### **Address for Correspondence**

**SUGUNA.T**

Research Scholar

M.Phil Mathematics

Bharath Institute Of Higher Education And Research

Mail Id: [rsugu1971@gmail.com](mailto:rsugu1971@gmail.com)

Guide Name: **Dr. R. DEEPA**

Assistant Professor, Department Of Mathematics

Bharath Institute Of Higher Education And Research

### **ABSTRACT**

Discrete outcome variable is common in public health, behavioral sciences and in many medical applications; the Poisson regression model is useful to analyze discrete random variable. For clustered discrete outcome variable where the observations are correlated among individual subjects, the number of observed discrete is sometimes greater than the expected frequency of the Poisson distribution and the discrete random variables are over-dispersed. Overdispersion is familiar in discrete random variable models particularly within the area of ecology and biological science because of missing covariates, non-independent, aggregations of data and an excess frequency of zeros. Every cluster levels received a singular level of a random effect that models the extra Poisson variation given within the data, are usually utilized to discuss heterogeneity in discrete random variable. However, studies investigating that the power of cluster level random effects as a way to discrete random variable model with over-dispersion is scarce. A situation where the variance of the response variable exceeds the mean, and hence, both over- dispersion and heterogeneity problems occur, in the appropriate imposition of the Poisson model may underestimate the standard error and overestimate the significance of the regression parameters, and so, giving misleading inference about the regression

*Research Paper*

parameters. Often, because of the hierarchal study design procedure, zero- inflation, over-dispersed and lack of independence may occur simultaneously, which render the standard ZIP model is inadequate. The multilevel ZIP and MZINB regression model are suitable for examining clustered correlated and over- dispersed discrete random variable with many zeros. In this thesis, derivation of a proposed score test for evaluating the over-dispersed, heterogeneity and zero-inflation parameters in discrete random variable regression models are performed.

**INTRODUCTION****BACKGROUND OF THE STUDY**

Modeling is the heart of applied Mathematical and Statistical sciences. Model is the key component in any Mathematical and Statistical analysis. In recent years in all most all fields of science, several research works have been directed to either the Mathematical models or the Statistical models. Model is a set of structural and functional relationships that can be expressed in terms of mathematical equations. A Mathematical model is a set of mathematical equations concerns with two or more variables. By introducing an error random variable or a disturbance term, the mathematical model becomes a stochastic model or statistical model.

Generally, the mathematical model or the statistical model may be specified either in the form of a linear model or in the form of a nonlinear model. The linear model consists of a set of linear equations concerns with two or more variables. Linear model has received greatest attention both in theory and in practice. From the theoretical point of view, it is mathematically tractable, and in practical applications of the wide variety, it has shown itself to be of great value. Further, many non-linear models can often be rearranged to be in a linear form.

Regression method is a statistical technique for investigating and modeling the relationship between the dependent and independent variables. Linear regression analysis establishes an average linear relationship between a dependent variable and a set of independent variables. Applications of linear regression models have a vital role in

analyzing various mathematical and statistical problems on different fields of science such as Economics, Business, Management, Engineering, Agriculture, Medicine, Social science, Biological and Life sciences, Physical sciences and Technology.

**STATEMENT OF THE PROBLEMS**

Nested data are very common in social sciences, psychology, health science, and others. The classical

*Research Paper*

linear models assuming that the observations have the same effects across groups. When the observations have different effects across the groups, the linear regression model is not appropriate. Ignoring group membership and focuses exclusively on inter-individual variation and on individual level attributes. This approach has a drawback of ignoring the potential importance of a group-level attribute in influencing an individual-level outcome. In addition, if outcomes for individuals within groups are correlated, the assumption of independence of observation is violated, resulting in an incorrect standard error and inefficient estimate (Diggle PJ et al., 1994).

The method focuses exclusively on inter-group variation and data aggregated to the group level variables. This approach eliminates the non-independence problem mentioned above but has the drawback of ignoring the role of individual-level variables

in shaping the outcome. Both methods essentially collapse all variables to the same level and ignore the multilevel structure. These approaches allow and define separate regression for each group coefficients to differ from group to group, but does not examine how specific group-level properties may affect individual-level outcomes or interact with individual-level variables.

**REVIEW OF LITERATURE**

The purpose of this chapter is to give a brief summary of concepts and theories about the generalized linear model for clustered discrete random variable, overdispersion, zero-inflation and the power of tests. This chapter also explains the review of the literature on the discrete random variable regression models.

**LINEAR STATISTICAL MODELLING**

Regression analysis is collection of statistical techniques for modeling and investigating the relationship between a response variable of interest and a set of predictor variables. In the classical regression model, the response variable  $y$  which is our main interest, select a sample of size  $n$  from our population of interest and observe values  $y_i$ ,  $i = 1, \dots, n$ , then wish to infer properties of the variable  $y$  in terms of other observed predictors  $x_i = (x_{1i}, \dots, x_{ki})$ . The main use of these predictor variables is to account for differences in the response variable or to put it another way to explain the variation in  $y$ . Consider the classical linear model

$$y_{ij} = X_{ij}\beta + e_i \dots\dots\dots (2.1)$$

where  $\beta$  is the coefficient of regression for  $X_{ij}$ , and  $e_i$  is a random or error term and assumed that the error terms are identically and independently distributed a normal distribution

*Research Paper*

with mean zero and variance  $\sigma^2$ . Also assumed that the outcome variable  $y$

$e$

$ij$

follows a normal distribution with mean  $X_{ij}\beta$  and variance  $\sigma^2$ . The coefficient of determination

( $R^2$ ) is a measure of the amount of variance in the dependent variable explained by the independent variable(s).

## **SCORE TEST FOR HOMOGENIETY OF GROUPS IN THE MULTILEVEL POISSON MODEL FOR CLUSTERED DISCRETE RANDOM VARIABLE**

### **3..1. INTRODUCTION**

Data with a multilevel nature often happens in public health, health service research, behavioral sciences, and in medicine. Examples embrace patients nested within hospitals, residents nested within a geographical area, students nested within a class, class nested within facilities and staff nested within a corporation. For clustered discrete random variable, the observations are usually correlated, there are three main approaches for correlated discrete random variable have been proposed (Perntile, 1988); conditional models, random effects model, generalized estimating equations (GEE`s). Conditional models Rosner (1984) are convenient only for particular cases, such as data with small group sizes or with order structure.

Analysts are progressively aware that the multilevel regression models are an appropriate way to analyze clustered data. A consequence of clustering of the groups within clusters or higher level units is that subjects from the same area could have more similar outcomes than subjects came from a different area. The multilevel regression models incorporate cluster-specific random effects that account for the dependency of the observation by partitioning the total individual variance into variation as a result of the cluster. The likelihood ratio test between the ordinary regression model and mixed effects models may be used as a homogeneity test (Self and Liang 1987). This chapter focuses on a multilevel Poisson regression model approach, the distribution of the response variables is modelled conditionally in a group-specific parameter that is itself a random variable (Laird and Ware 1982, Stiratelli, Laird, and Ware 1984), to take of the coefficient of regression and random parameters in Poisson discrete points.

### **THE MULTILEVEL POISSON REGRESSION MODEL**

In the model formulation, when the outcome variable is discrete denoting the number of time that an incident occurred, a Poisson regression model can be accustomed relate the mean number of events to

*Research Paper*

a group of explanatory variables employing a logarithmic link function. Then, the Poisson regression model will be used as  $lo(\mu_i) = X_i\beta, Y_i \sim Poisson(\mu_i)$  ----- (3.1)

where  $X_i$  denotes a  $p \times 1$  column matrix of covariates measured with the  $i^{th}$  subject,  $Y_i$

denotes the discrete outcome variable measured with the  $i^{th}$  subject, denotes a  $1 \times p$  row matrix of the regression coefficients and the parameter  $\mu_i$  denotes the expected or mean number of events for the  $i^{th}$  subject given their set of observed covariates. Consider a two-level random intercept Poisson regression model and it is assumed that the intercept is allowed to vary randomly across the groups.

Let  $Y_{ij}$  be the response variable for the observation  $j$  of group  $i, i = 1, 2, \dots, k; j = 1, 2, \dots, n_i$

$k$

with  $N = \sum_{i=1}^k n_i$ , the conditional distributions of the outcome variable  $i = 1$

$y = (y_1, y_2, \dots, y_{i1}, y_{i2}, \dots, y_{in_i})$ , given a set of cluster level random effects  $u_i$  and the probability

density functions of  $Y_{ij}$  is defined as follows

$$f(Y_{ij} | \mu_{ij}) = \prod_{j=1}^{n_i} \frac{\mu_{ij}^{y_{ij}} \exp(-\mu_{ij})}{y_{ij}!} \quad \mu_{ij} = \exp(\beta + z_{ij} u_i)$$

$n_i$

$$= \exp \left\langle \sum_{j=1}^{n_i} \{y_{ij}(x_{ij}\beta + z_{ij}u_i) - \exp(x_{ij}\beta + z_{ij}u_i) - \log(y_{ij}!)\} \right\rangle$$

$i=1$

$n_i$

$$= \exp \left\langle \sum_{j=1}^{n_i} \{\theta_{ij} y_{ij} - g(\theta_{ij})\} \right\rangle + C(y_{ij}, c_i) = \dots \dots \dots (3.2)$$

*Research Paper*

i=1

*Research Paper*

where  $(y_{ij}, c_i) = \sum_{i=1}^{n_i} \log(y_{ij}!)$ , which  $C(y_{ij}, c_i)$  does not depend on the model parameters. The mean and the variance of  $Y_{ij}$  are respectively,

$$\mu_{ij} = E(y_{ij}/\alpha_i) = g'(\theta_{ij}) = \exp(-x_{ij}\beta + z_{ij}\alpha_i) \text{ and}$$

$$\sigma_{ij}^2 = \text{Var}(y_{ij}) = c_i g''(\theta_{ij}) = \mu_{ij} = \exp(-x_{ij}\beta + z_{ij}\alpha_i) /$$

where 'denotes the differentiation with respect to parameter  $\theta_{ij}$ . Consider the mixed effects model allowing at least one regression coefficient to be random is

$$\theta_{ij} = \log(\mu_{ij}) = \log(\alpha_i) + x_{ij}\beta + z_{ij}\alpha_i \dots \dots \dots (3.3)$$

where  $\beta$  denotes a  $p \times 1$  vector of fixed effects with its associated design vector  $x_{ij}$  and

$\alpha_i$  is the scalar random subject with associated covariates  $z_{ij}$ . In this model the mean

$\mu_{ij}$  and the variance  $\sigma_{ij}^2$  are conditionally on  $\alpha_i$ , now to test homogeneity across and within groups, consider the random intercept model in which  $z_{ij}=1$  for all  $i, j$ . To test

the homogeneity hypothesis that all of the variables and the correlation among the random effects are zero in a generalized linear mixed model.

The parameter  $\alpha_i$  can be written as  $\alpha_i = \alpha + D^{1/2}u_i$ , where the  $u_i$ 's are independently and identically distributed as a normal distribution with zero mean and unit variance. Therefore,  $\alpha_i$ 's are identical and independently distributed with mean  $\alpha$  and variance  $D$ . Our interest is to test  $H_0 : D = 0$  against the alternative  $H_0 :$

$D > 0$ . Note that for discrete random variable models; this is equivalent for testing homogeneity across groups as well as testing homogeneity within groups.

The first and the second partial derivatives of the log likelihood function with respect to the scalar random subject parameter  $\alpha_i$  for the multilevel Poisson regression model is given by

Research Paper

$$\frac{\partial}{\partial \alpha_i} \left[ -\log f_{ij}(\beta, \alpha, D=0) \right] = \sum_{j=1}^{n_i} (y_{ij} - \mu_{ij})$$

and

$$\frac{\partial^2}{\partial \alpha_i^2} \left[ -\log f_{ij}(\beta, \alpha, D=0) \right] = - \sum_{j=1}^{n_i} \frac{1}{\sigma_{ij}^2}$$

Using the first and second partial derivation of the log likelihood equation with respect to  $\alpha_i$ , the score statistic is given by

$$S_{\beta, \alpha} = \sum_{i=1}^k \left\{ \sum_{j=1}^{n_i} (y_{ij} - \mu_{ij}) z_{ij} - \sum_{j=1}^{n_i} [z_{ij} \sigma_{ij}] \right\}$$

Then the score test statistic for testing homogeneity  $H_0: D = 0$  for the known nuisance parameters  $\beta$  will be

$$H_{PD} = \frac{S_{PD}^2(\beta, \alpha)}{I_{D\beta} I_{\beta\beta}^{-1} I_{\beta D}} \dots \dots \dots (3.5) \quad (I_{DD} - I_{D\beta} I_{\beta\beta}^{-1} I_{\beta D})$$

Now, the asymptotic variance function as the group size  $k \rightarrow \infty$  of  $S_{PD}(\beta, \alpha, D)$  under  $H_0$  (Cox and Hinkley, 1974) can be expressed as a function of information matrix. Then the asymptotic variance function for the score test is

$$I(\beta) = I_{DD} - I_{D\beta} I_{\beta\beta}^{-1} I_{\beta D}$$



*Research Paper*

where  $I = \sum_{i=1}^k \frac{\partial^2 l}{\partial \beta_i \partial \beta_i} \Big|_{D=0}$  is a scalar

$$I_{\beta\beta} = \sum_{i=1}^k \frac{\partial^2 l}{\partial \beta_i \partial \beta_i}, \quad I_{\beta D} = I_{D\beta} = \sum_{i=1}^k \frac{\partial^2 l}{\partial \beta_i \partial D}$$

**PARAMETRIC ESTIMATIONS OF SCORE TEST BASED ON THE MULTILEVEL POISSON REGRESSION MODEL**

Now to evaluate the variance of the score function (S) defined as  $Var(D) = I_{DD} - I_{D\beta}I_{\beta D}$ . The  $i^{th}$  summand of  $I_{DD}$  can be written as

$$E \left( \frac{\partial^2 l}{\partial D^2} \mid D = 0^* \right) = \sum_{i=1}^2 E \left\{ \sum_{j=1}^{n_i} (y_{ij} - \mu_{ij})^2 - \sum_{j=1}^{n_i} \sigma_{ij}^2 \right\}$$

To solve the variance of the score function, let us to define

$U_{ij} = y_{ij} - \mu_{ij}$ , the  $i^{th}$  term of the score test can be written as

$$\frac{\partial l_i}{\partial D} = \sum_{j=1}^{n_i} U_{ij} - \sum_{j=1}^{n_i} \mu_{ij}$$

Thus

$$E \left( \frac{\partial^2 l_i}{\partial D^2} \mid D = 0^* \right) = \sum_{j=1}^{n_i} \sigma_{ij}^2$$

*Research Paper*

where  $U_i = \sum_{j=1}^{n_i} U_{ij}$ , now since  $(U_{ij}) = 0$ , then

$$E(U^2) = E \left( \sum_{i=1}^{n_i} U_i^2 \right) = E \left( \sum_{i=1}^{n_i} \sum_{j=1}^{n_i} U_{ij}^2 + \sum_{i=1}^{n_i} \sum_{j=1}^{n_i} U_{ij} U_{ij'} \right) = E \left( \sum_{i=1}^{n_i} \sum_{j=1}^{n_i} U_{ij}^2 \right)$$

$$= E \left( \sum_{i=1}^{n_i} \sum_{j=1}^{n_i} (U_{ij} - \mu_{ij})^2 + \sum_{i=1}^{n_i} \sum_{j=1}^{n_i} (U_{ij} - \mu_{ij})(U_{ij} + \mu_{ij}) \right) = \sum_{i=1}^{n_i} \sum_{j=1}^{n_i} \sigma_{ij}^2$$

Therefore,

$$E \left( \frac{\partial l}{\partial D_4} \right) \Big|_{D=0^*} = \frac{1}{2} E^* U^2 - (U^2)_i^2 = \frac{1}{4} Va(U^2) = \frac{1}{4} (\mu_4 - \mu^2)$$

where  $\mu_2$  and  $\mu_4$  are the second and the fourth central moments of  $Y_{ij}$ , respectively, which can be expressed as a function of the second and the fourth cumulates  $K_2$  and

*Research Paper*

$K_4$  of  $Y_{ij}$  (Kendall and Stuart, 1977).  $\mu_4 = K_4 + 3K_2^2 = K_2 = \sigma^2 = \mu_{ij}$ , then

$\mu_4 = \mu_{ij} + 3\mu^2$ . After simplification the values of  $I_{DD}$

$$\frac{1}{4} \sum_{i=1}^k \sum_{j=1}^{n_i} (\mu_{ij} + 2\mu_{ij}^2)$$

The variance of the score functions can be derived from the Fisher information.

$$I(\beta) = \begin{pmatrix} I_{\beta\beta} & I_{\beta D} \\ I_{D\beta} & I_{DD} \end{pmatrix} *$$

where

$$\frac{\partial^2 l}{\partial \beta^2} = \sum_{j=1}^{n_i} \frac{\partial^2 l_{ij}}{\partial \beta^2} = \sum_{j=1}^{n_i} \left( -\sum_{i=1}^k \mu_{ij} x_{ij} \right)$$

Therefore

$$k \quad -\hat{\alpha}$$

Research Paper

$k$   $n_i$

$$-\partial^2 \log f(\beta, \alpha, D) = 0$$

$$I_{\beta\beta} = \sum_{i=1}^k E \left[ \frac{\partial^2 \log f}{\partial \beta \partial \beta'} \right] = \sum_{i=1}^k \sum_{j=1}^k E \left[ \frac{\partial^2 \log f}{\partial \beta \partial \beta'} \right]$$

$$= \sum_{i=1}^k \sum_{j=1}^k g''(\theta_{ij}) x_{ij} x_{ij}' = \sum_{i=1}^k \sum_{j=1}^k \mu_{ij} x_{ij} x_{ij}'$$

$$I_{\beta D} = \sum_{i=1}^k E \left[ \frac{\partial^2 \log f}{\partial \beta \partial D} \right] = \sum_{i=1}^k E \left[ \frac{\partial^2 \log f}{\partial \beta \partial D} \right] = \sum_{i=1}^k \sum_{j=1}^k y_{ij} x_{ij} - \sum_{i=1}^k \sum_{j=1}^k \mu_{ij} x_{ij}$$

$$I_{\beta D} = \sum_{i=1}^k \sum_{j=1}^k g'''(\theta_{ij}) x_{ij} = \sum_{i=1}^k \sum_{j=1}^k \mu_{ij} x_{ij}$$

where  $x_{ij} = (1, x_{1ij}, \dots, x_{n_{ij}})$ .

which asymptotically as  $n \rightarrow \infty$  has a chi-square distribution with one degree of freedom, now the maximum likelihood estimate of  $\beta$  can be estimated iteratively by using Fisher's scoring method from the following equations.

**MODEL SELECTION**

If there are several models to be compared, in order to select the best model which fits the data instead of using the likelihood ratio test, it can be selected by using the Akaike information criteria (AIC) and Bayesian information criteria (BIC).

**3.7.1 AKAIKE INFORMATION CRITERIA**

AIC is the most common means of identifying the model which fits the data well by comparing two or more than two models. The goodness of fit test against the complexity of the model is similar to that of the coefficient of multiple determination ( $R^2$ ); however, it is penalized by the number of parameters included in the complexity of the model. Unlike the  $R^2$ , the good model is the one which has the minimum AIC value. It is given by the following formula  $AIC = -2\ell + 2k$ , where  $\ell$  is the log likelihood function of a model that will compare with the other models and  $k$  is the number of parameters in the model including the intercept (Ismail and Jemain, 2007).

**BAYESIAN INFORMATION CRITERIA**

Unlike the Akaike information criteria, the Bayesian information criteria take into account the size of the data under consideration. It is given by  $BIC = -2\ell + k \log(n)$  where  $\ell$  is the log-likelihood of a model that will compare with the other models,  $n$  is the sample size of the data and  $k$  is the number of parameters in the model including the intercept.

**SIMULATION STUDY**

In this section, a simulation study is conducted to compare the proposed and the existing models in terms of sizes and powers. For studying the properties of the statistic in terms of empirical size, generating discrete random variable from a Poisson distribution under the null hypothesis of homogeneity and assume that the random effects parameter are one ( $z_{ij} = 1$ ) and the samples are comprised of 10; 20; 50;

100 observations and 5; 10; 20; 50 groups with simulating data, the multilevel Poisson distribution is simulated for the distribution of the response variable assuming that under the null hypothesis homogeneity within different number of groups according to the variance  $D$  of the distribution of the group-specific random effects of the response variable. Each simulation experiment for level and power was based on 1000 simulated samples.

In the simulated model, we reviewed the literature to define initial parameter values. We presented the result of a small simulation study examining the empirical size and power of the test statistics discussed in this thesis. The following log linear model for the response variable is assumed (see

*Research Paper*

Jacqmin-Gadda and Commenges (1995)).  $\log(\mu_{ij}) = 0.8x_{1ij} + 0.5u_i - 0.5$  (3.8)

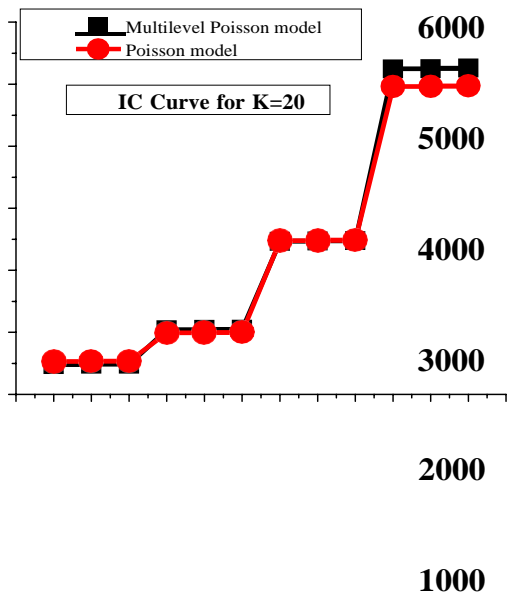
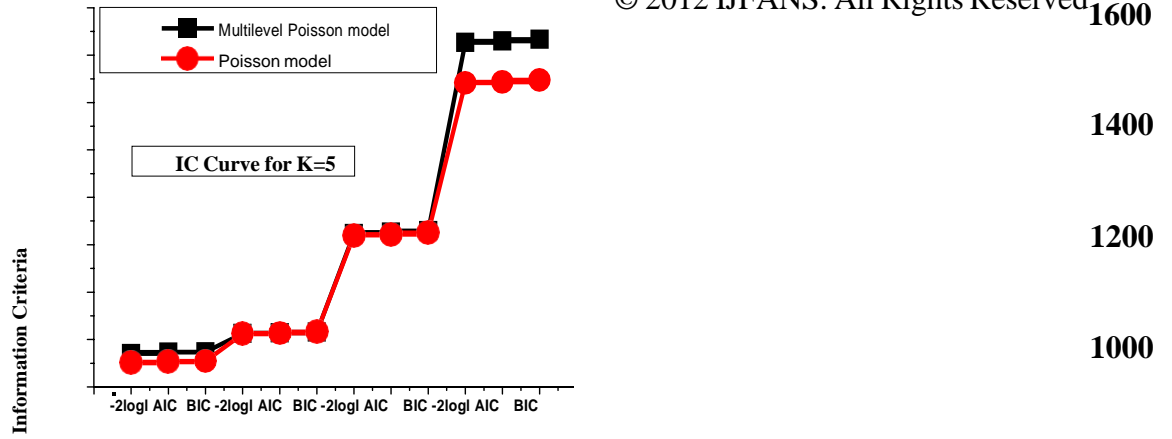
$y_{ij} \sim \text{Poisson}(\mu_{ij})$ ,  $i = 1, 2, \dots, k$ ;  $j = 1, 2, \dots, n_i$ ,

For the variable  $x_{1ij}$  is a subject-specific effect which is simulated from a uniform distribution with mean zero and unit variance and  $u_i$  is group specific effects and simulated from a standard normal distribution and generate a set of random numbers from a uniform distribution in the interval (0,1) as the values of  $x_{1ij}$ . For drawing samples and for estimating the maximum likelihood estimates of the regression and homogeneity parameters of interest under the null hypothesis,  $a_i = a + D^{1/2}u_i$ .

To simulate correlated data, added a group-specific random effect under the hypothesis of homogeneity, that is,  $D = 0$ , where the  $u_i$ 's are identical and independently distributed with a standard normal distribution with mean zero and unit variance. Therefore,  $a_i$ 's are identical and independently distributed with mean  $\alpha$  and variance  $D$ . For each set of generated data, a multilevel Poisson model is fitted for calculating the score test and the existing tests followed by the powers of the score tests. Results from the simulation study are presented in Table- 3.1 and 3.2.

In Table 3.1 the results investigated that how the information criteria perform in the multilevel Poisson regression model selection problems via simulations. Model with smaller AIC is considered to be better. When the number of clusters is large and the number observation is small then the multilevel Poisson regression model is better than the standard Poisson regression model, whereas when the sample size and cluster numbers are small then the Poisson regression model is better than the multilevel Poisson regression model. The simulated data results indicated that the performance of the criteria to select the true model generally involved with an increase of sample size, despite differences in performance among the information criteria. The simulated results indicate that the performance of the model depends on the sample size and the number of clusters.

Research Paper



10 20

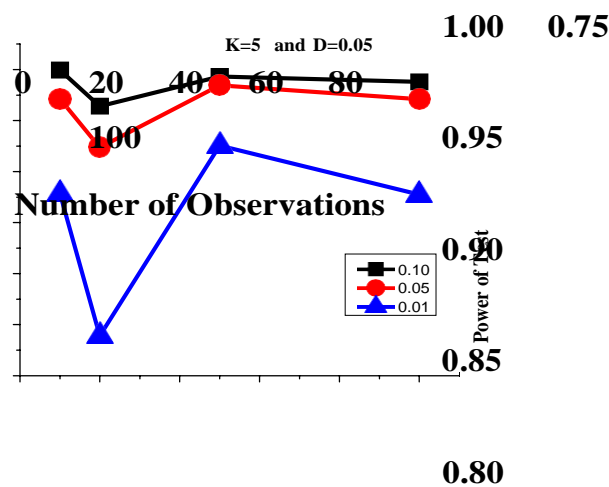


*Research Paper*

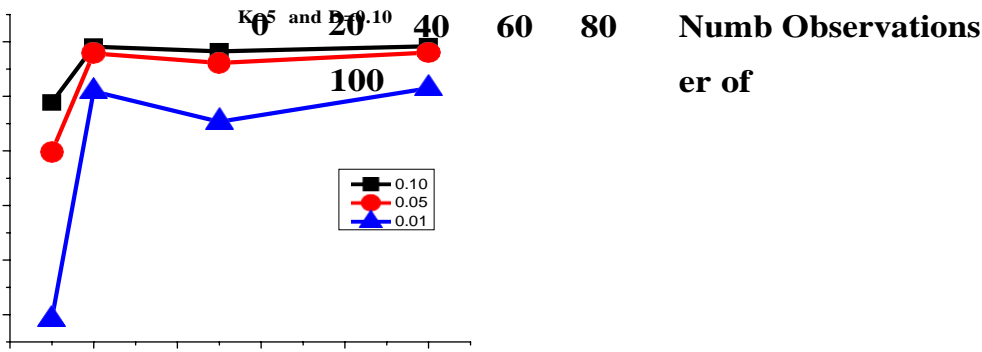
Data are generated from the multilevel Poisson distributio under the nullhypothesis on 1000 replication. In the simulation, thelevels (10%; 5% and 1%), the sample size (10, 20, 50 and 100) and the number of groups (5, 10, 20 and 50 ) areconsidered.

In note that an error probability increases power increases. The power of the tests of the three scenario increases when a increases, and for large sample groups and small variance for the group effect ( k = 50, n = 10, D=0.05) the power increase fast and approaches to 1. For small sample groups, when the standard deviations of the group effectsincrease, the power increases slowly, whereas in large sample groups, (k = 50, n = 50, D=0.15 and a =0.1) as standard deviations of the group effect increases, the power increases slowly, however, when the values of D increase from 0.05 to 0.15, the power decreases. Generally, as the number of groups increases, the power is slightly increases. Therefore, the proposed score test is more important for testing and controlling heterogeneity of the group effects by fixing the number of observation and number of groups due to its high power to predict the model.

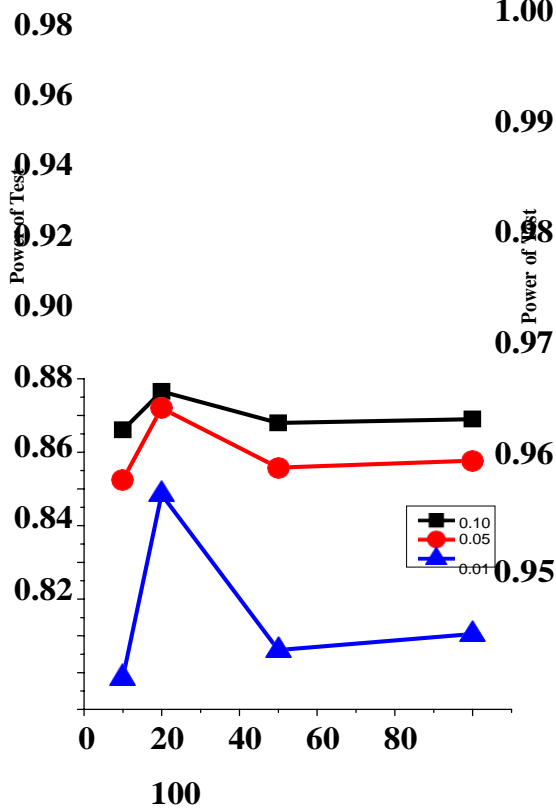
Increasing the sample size will decrease the standard error (and increasepower). Similarly, increase the amount of variance in x will increase power. However, increasing amount of unexplained variance will serve to decrease power. Sample size is not the only factor for power in the multilevel model, effects on power can depend on the parameters (intercept estimate, individual specific intercept, slope estimate, individual specific slope, standard error of the average (between person), slope error, variances of the independent variable, variances of individual slope (multilevel), sample size and cluster size), Sean P. Lane and Erin P. Hennes (2018).



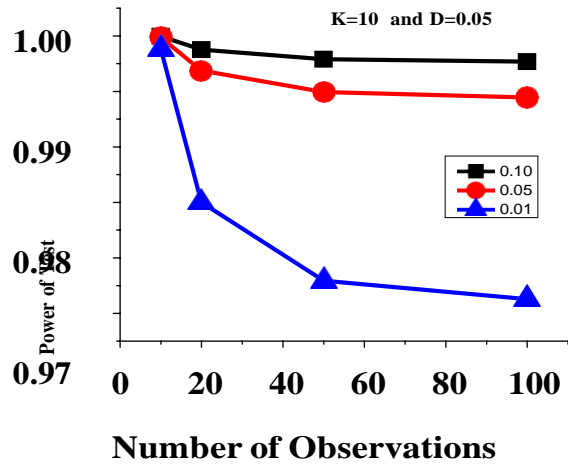
Research Paper



1.00 K=5 and D=0.15

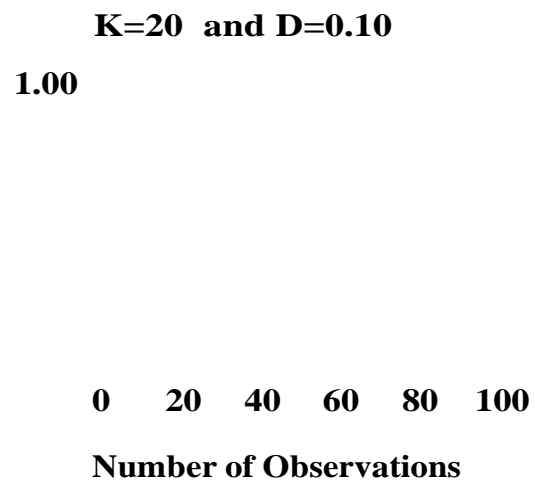
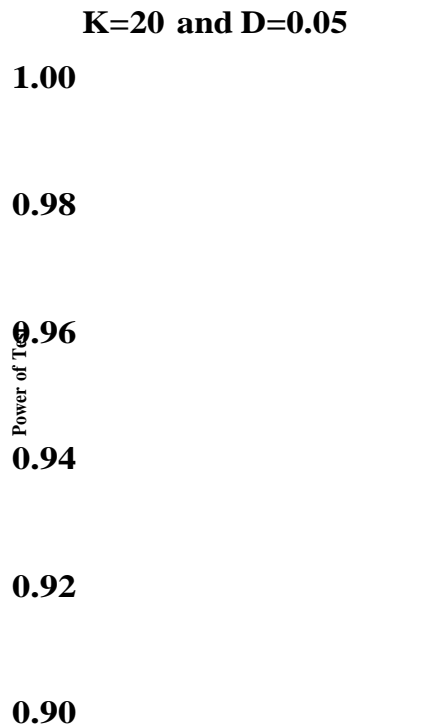
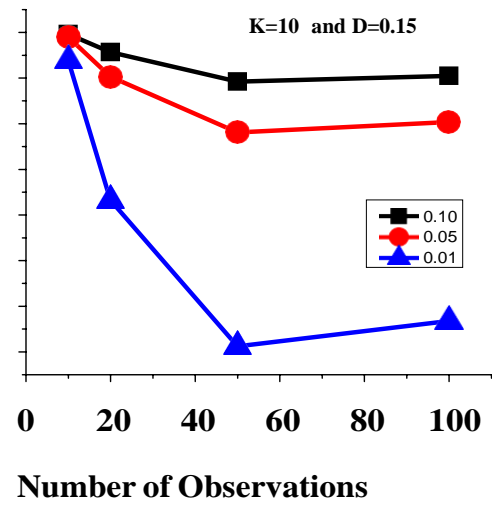
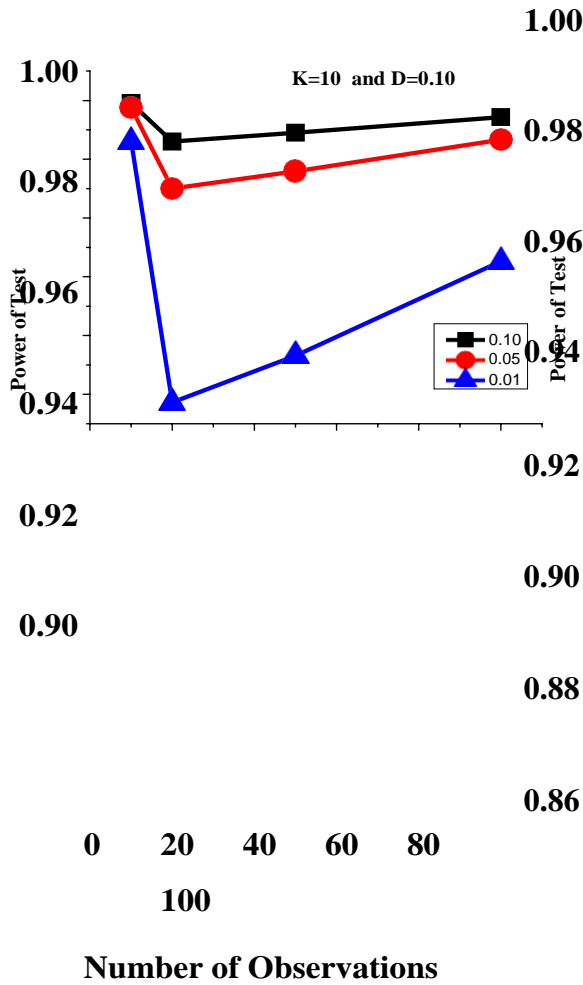


Number of Observations



Number of Observations

Research Paper



Research Paper

0.98

0.96

0.94

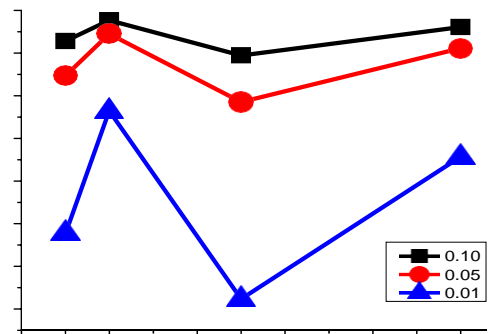
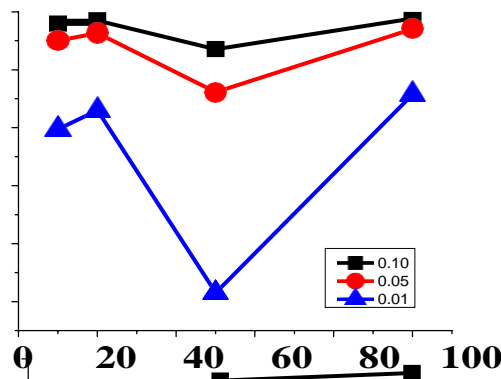
0.92

0.90

0.88

0.86

Power of Test



Number of Observations

Number of Observations

K=20 and D=0.15

1.00

0.95

0.90

0.85

0.80

0.75

0

0

0

0

0

00

2

4

6

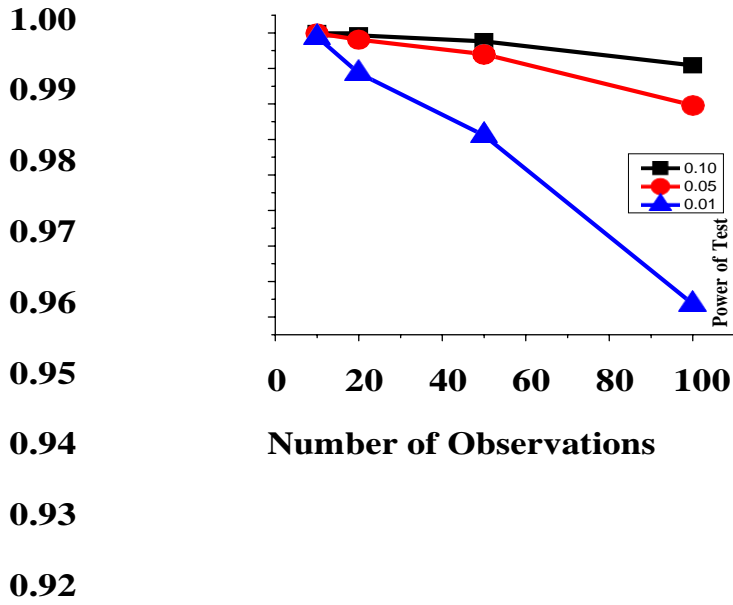
8

1

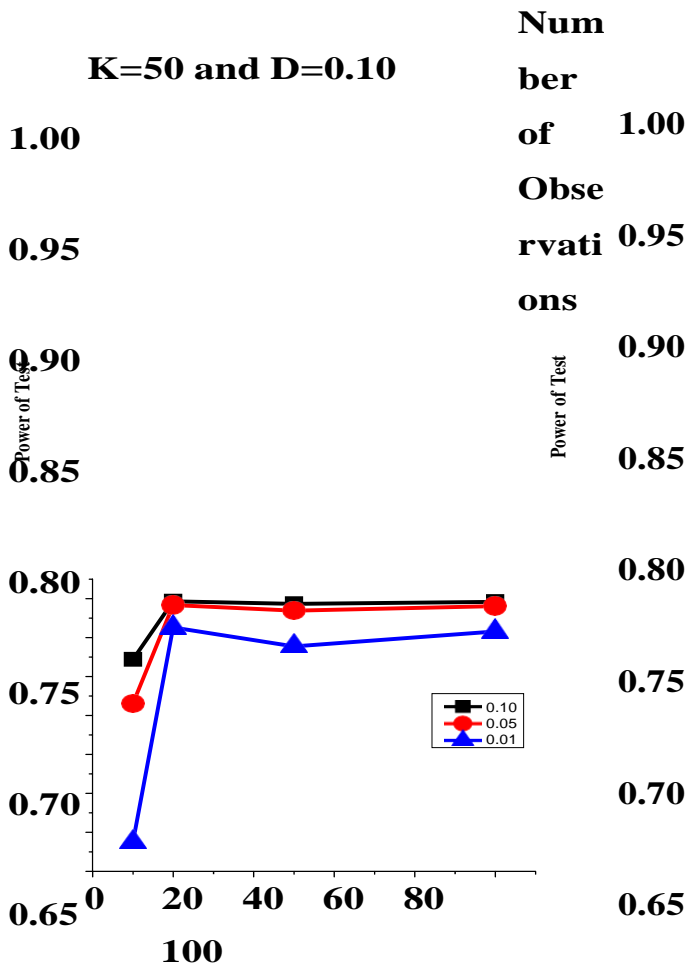
00

Research Paper

**K=50 and D=0.05**

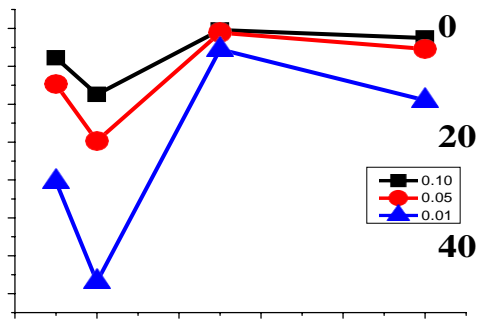


**K=50 and D=0.10**



*Research Paper*

**K=50 and  
D=0.15**



**60 80 100**

**Number of Observations**

**Figure 3.2. Simulated curve for power and goodness of fit test of the models**

*Research Paper*

In this article, the power of the multilevel Poisson regression models via simulated data results are presented. A simulation and application data are used to illustrate our method. The results revealed that the proposed score test is more preferable than the existing model. Based on the information criteria, AIC, the

proposed model is better than the standard Poisson model. The results in Table 3.8, revealed that deaths of children varied among regions. In addition, all covariates and dummy explanatory variables were found to be significant difference in the deaths of children between regions. The AIC values of the empty model with random intercept are larger than that of the random intercept and fixed coefficient model, which implies that the random intercept model is better than the standard Poisson model and also on the predicted probability, the multilevel Poisson regression model is better than a Poisson regression model.

**CONCLUSION**

In clustered discrete random variable, when the responses of each observation are correlated, familiar ANOVA and regression models do not give optimal analysis. The standard multilevel models yield correct inference for clustered normally distributed data. Generalized linear models specifically the Poisson and negative binomial regression models give correct inference for non clustered data. In this thesis, developed the multilevel discrete random variable models illustrated with examples and simulation study, presenting some score test statistics for analyzing overdispersion, zero-inflated and heterogeneity in clustered discrete random variable and compare the existing model with its alternative model tests and identified the best test statistic for testing coefficients of regression, overdispersion, zero inflation, and heterogeneity parameters of clustered discrete random variable in terms of its power and size.

We develop a proposed score test based on the multilevel Poisson model for testing heterogeneous parameter in the equidispersed clustered discrete random variable and analyzed the Poisson regression model with the assumption that it is used for model fit under the null hypothesis. Furthermore, the likelihood ratio tests are used as an alternative test to select the best model. From the simulation and application study results shown that when the dataset has heterogeneous groups in discrete random variable, the multilevel Poisson regression model gives a good and correct result in the analysis while Poisson regression is clearly not adequate for handling heterogeneous data. However, when

*Research Paper*

the data has homogenous between groups, the Poisson regression model is more reliable. From the simulation study, for fixed values of sample size, when the coefficient of regression and heterogeneity parameter are increasing the power of the scores are increasing. On the other hand, for fixed values of the coefficient of regression and heterogeneous parameters, when the sample size increasing the power of the score test is increasing. For large values of the sample size and coefficient of the regression and heterogenous parameters, then the difference among different tests become trivial in terms of its power. For other cases, the proposed score test is more appropriate for general use because of its high Power.

**REFERENCES**

- Abramowitz, M., and Stegun, I.A. (1972). Handbook of Mathematical Functions, New York: Dover Publications, Inc.
- Agresti, A. (1996). An introduction to categorical data analysis, John Wiley and Sons, New York.
- Alker, H. R. (1969). A typology of ecological fallacies, in Dogan, M. and Rokkan, S. (eds.), Quantitative Ecological Analysis in the Social Sciences. Cambridge, MA: MIT Press.
- Anseombe, F.J. (1967). Topics in the Investigation of Linear Relations Fitted by the Method of Least Squares, *J.R.S.S., Series*, 29, 1-32.
- Bohning, D. (1998). Zero inflated models and C.A.MAN: a tutorial collection of evidence, *Biometrical*, 40, 833-843.
- Bolger, N., Zuckerman, A., & Kessler, R. C. (2000). Invisible support and adjustment to stress. *Journal of Personality and Social Psychology*, 79, 953– 961. Sage Publications, New York
- Bolker BM, Brooks ME, Clark CJ, Geanges SW, Poulsens MHH, White JSS. (2009). The generalized linear mixed models: a practical guide for ecology and evolution. *Trends in Ecology & Evolution*, 24, 127-135.
- Breslow, N. E. (1999). Tests of hypotheses in over dispersed Poisson regression and other quasi-likelihood models, *Journal of the American Statistical Association*, Vol. 85, 565-571.
- Bryk, A.S., & Raudenbush, S. W. (1992). Hierarchical linear models (applications and data analysis methods). Sage Publications, New York.
- Cameron, A. Colin & Trivedi, Pravin K., (1990). Regression-based tests for overdispersion in the Poisson model. *Journal of Econometrics, Elsevier*, vol.46(3), 347-364.
- Diggle, P.J., Liang, K.Y. and Zeger, S.L. (1994). The Analysis of Longitudinal Data.



*Research Paper*

Oxford. Clarendon.

- Draper, N.R. and Smith. If. (1998). Applied Regression Analysis. Third Edition, John Wiley & Sons. New York.
- Efronson. M.A. (1960). Multiple Regression Analysis. In A.Ralston and H.S. Wile (Eds.) Mathematical Methods for Digital Computers, John Wiley, New York
- Fisher, R. A. (1941). Average excess and average effect of a gene substitution. *Annals of Eugenics* 11, 53–63.
- Fitzmaurice GM, Laird NM, Ware JH(2004). Applied longitudinal analysis. New Jersey: John Wiley & Sons Inc
- Gardner, W., Mulvey, E.P., and Shaw, E.C (1995). Regression Analyses of Discreted and Rates: Poisson, Over-dispersed Poisson, and Negative Binomial Models, *Psychological Bulletin*, 118, 392-404.
- Hilbe JM.(2011). Negative binomial regression (2<sup>nd</sup> edition). Cambridge: Cambridge University Press.
- Hinde, J. P. and Demetrio, C. G. B. (1998). Overdispersion: models and estimation, *Computational Statistics and Data Analysis*, Vol. 27, 151-170.
- Hocking, R.R. (1976). The Analysis and Selection of Variables in Linear Regression, *Biometrics*, 32, 1-49.
- Jacqmin-Gadda, H., and Commenges, D.(1995). Tests of homogeneity for generalized linear models. *Journal of the American Statistical Association*, 90, 1237 – 1246.
- Jansaku, N. and Hinde, J.P.(2004). Linear mean variance negative binomial models for analysis orange tissue culture. *Journal of Science and Technology*, Vol. 26, No.5, 683 – 689.
- Jansaku N. and Hinde, J.P.(2002). Score Tests for Zero Inflated Poisson Models, *Comput. Data Anal*, 40, 75-96.
- Jansaku, N., and Hinde, J.P.(2009). Score Tests for extra-Zero models in zero inflated negative binomial models, *commun. Statist. Simul. Comput*, 38, 92-108.
- Kendall, M.G. and Stuart, A. (1977). The Advanced theory of statistics. Macmillan, New York, Vol. 1, 168.
- Kibria, B.M.G., Kristofer, M. and Shukur, G. (2013). Some Ridge Regression Estimations for zero inflated Poisson Model. *Journal Applied Statistics*. 40(4), 721-735.
- Leckie, G. and Goldstein, H. (2009). The limitations of using school league tables to inform school choice, *Journal of the Royal Statistical Society: Series A*, 172, 835-851.

*Research Paper*

- Lee, A.H., Wang, K., Yau, K.K.W., Carrivick, P.J.W., and Stevenson, M.R. (2005). Modelling bivariate discrete series with excess zeros. *Mathematical Biosciences* 196, 226 –237.
- McCulloch, C. E., Searle, S. R. & Neuhaus, J. M. (2008). Generalized, Linear, and Mixed Models. Hoboken, New Jersey: John Wiley & Sons, Inc., 2nd ed.
- McGilchrist, C.A. (1994). Estimation in generalized mixed models, *J.Roy. Statist.Soc. Ser.B* 56, 61-69
- McLachlan G. J., and Krishnan, T. (1997). The EM Algorithm and extensions, Wiley, New York.
- Neyman, J. (1959). Optimal asymptotic tests for composite hypothesis, *In Probability and Statistics*, The Harold Cramer, 213-34, U. Grenander (ed), New York, Wiley.
- Neyman, J. and Scott, E. L. (1966). On the use of ( $\alpha$ ) Optimal tests of composite hypothesis, *In Bulletins of the Institute of the International Statistics*, Vol. 41, 477-497.
- Osgood, D. W., and Chambers, J. M. (2000). Social disorganization outside the metropolis: An analysis of rural youth violence. *Criminology* 38: 81–115.
- Paul, S. R. and Banerjee, T. (1998). Analysis of two-way layout of Discrete random variable involving multiple discretely in each cell. *Journal of the American Statistical Association*, Vol. 93, 1419-1429.
- PAUL, S.R., and AZAD, K. (2011). Testing Homogeneity in Clustered (Longitudinal) Discrete random variable Regression Model with Over-Dispersion department of Mathematics & Statistics, University of Windsor, Windsor, Canada.
- Pearson, K. (1914). Life and Letters and Labours of Francis Galton, University Press, Cambridge.
- Pettiquanti. M. and Neqaraja, N. (2017). Power analysis for negative binomial models with application to multiple Sclerosis clinical trials, *Journal of binomial Biopharm Statistics*, Vol. 22, No.2.
- Rabe-Hesketh, S, and Skrondal, A. (2005). Multilevel and Longitudinal Modelling using Stata, Stata Press, Stata Corp, College Station, Texas.
- Rabe-Hesketh, S., Skrondal, A. and Zheng, X. (2012). Generalized multilevel structural equation modeling. In Hoyle, R. (Ed.). *Handbook of Structural Equation Modelling*. Guilford Press, 512-531.
- Rao, C. R. (1947). Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation, *Proceedings of the Cambridge Philosophical*

- Saleh, A.K.Md.E., Arashi, M., and Kibria, B.M.G.(2019). Theory of Ridge Regression Estimation with Applications. Wiley, New York.
- Sean P. Lane and Erin P. Hennes (2018). Power struggles: Estimating sample size for multilevel relationships research. *Journal of Social and Personal Relationships*, Vol. 35(1) 7–31. Purdue University, USA
- Thall, P. F. (1992). Score tests in the two-way layout of discretely, *Communications in Statistics, Theory and Methods* Vol. 21, 3017-3036.
- Trivedi,P.K. (1990). Memorandum of Understanding: An Approach to Improving Public Enterprise Performance, International Management Publishers.
- Van den Broek J. (1995). A score test for zero-inflation in a Poisson distribution, *Biometrics*, 51, 738 –743.
- Verbeke, G., and Molenberghs, G. (2000). Linear mixed models for longitudinal data, Springer Series in Statistics, Springer-Verlag, New-York.
- Vonesh, E.F. and Chinchilli, V.M. (1997). Linear and Nonlinear Models for Analysis of Repeated Measurements. Marcel Dekker, New York.
- Wang ,K.,Yau, K.K.W., and Lee, A.H. (2002). A zero inflated Poisson mixed model to analyze diagnosis related groups with majority of same day hospital stays, *Comput. Meth. Programs Biomed*, 68, 195 – 203.
- Wang, P., Puterman, M. L., Cockburn, I. and Le, N. (1996). Mixed Poisson regression models with covariate dependent rates, *Journal of Biometrics*, Vol. 52, 381-400.
- Xiang, L., Lee, A.H., Yau,K.K.W., and McLachlan, G.J. (2006). A score test for zero inflation in correlated Discrete random variable, *Statist.Med*, 25, 1660-1671.
- Xiang, L., Lee, A.H., Yau, K.K.W., and McLachlan,G.J.(2007). A score test for overdispersion in zero inflated Poisson mixed regression model, *Statist.Med*, 26, 1608-1622.
- Xie, M., He, B., and Goh, T. N. (2000). Zero-inflated Poisson model in statistical process control, *Computational Statistics and Data Analysis*, Vol. 38, 191-201.
- Yang, J., Benyamin, B., McEvoy, B. P., Gordon, S., Henders, A. K., Nyholt, D. R., Madden P. A., Heath, A. C., Martin, N. G., Montgomery, G.W., Goddard, M. E., &Visscher, P. M. (2010). Common SNPs explain a large proportion of the heritability for human height. *Nature Genetics*, 42, 565–569.
- Yau,K.K.W., and Lee, A.H.(2001). Zero inflated Poisson regression with random effects to evaluation an occupation injury prevention programs, *Statist. Med*, 20, 2907-2920.

*Research Paper*

- Yau, K.K., Lee, A.H. and Ng, A.S.K. (2003). Finite mixture regression model with random effects: Application to neonatal hospital length of stay, *Computational Statistics & Data Analysis* 41, 359 – 366.
- Zhao, Y., James, H., Cheryl, L. (2009). Score tests for over dispersion in zero-inflated Poisson mixed models, *Epidemiology and Biostatistics*, Arnold School of Public Health, University of South Carolina, SC 29208, USA.
- Zhu, H. And Zhang, H. (2006). Generalized score test of homogeneity for mixed effects models. *The Annals of Statistics*, 34, 1545 – 1569.
- Zuur, A.F., Ieno, E.N., Walker, N.J., Saveliev, A.A., Smith, G.M. (2009). *Mixed effects models and extensions in ecology with R*. New York: Springer.