

# Com-Poisson Model with AR(1)-type Correlation Structure for Longitudinal Count Response Data

Jowaheer, Vandna

*University of Mauritius, Department of Mathematics*

*Reduit, Mauritius*

*E-mail: vandnaj@uom.ac.mu*

Mamodekhan, Naushad Ali

*University of Mauritius, Department of Economics and Statistics*

*Reduit, Mauritius*

*E-mail: n.mamodekhan@uom.ac.mu*

## ABSTRACT

Analysis of longitudinal count responses, affected by one or more explanatory variables, requires adequate modeling of the correlations underlying the responses repeatedly collected over time. These correlation structures may be then used with the corresponding means and variances of the process to formulate quasi-likelihood estimating equations which can be solved to obtain reliable estimates of the regression parameters. True Gaussian ARMA-type correlation structures have been found to be more efficient against 'working' as well as 'random effects based' correlation structures in analyzing Poisson longitudinal counts. Quite often, count responses are under-dispersed or over-dispersed relative to Poisson distribution. Com-Poisson distribution, due to its convincing properties, has been widely used to describe under- or over- dispersed count data in cross-sectional set up. However, there exists no application of Com-Poisson model in longitudinal case. The challenge lies in modeling the correlation structure of repeated Com-Poisson counts. In this paper, we provide the framework for developing and analysing a Com-Poisson model with AR(1)-type correlation structure under the longitudinal set-up.

## INTRODUCTION

Com-Poisson distribution<sup>[10]</sup> is a weighted Poisson distribution, capable of modeling the counts subject to over-, under- or equi- dispersion. This capability of Com-Poisson distribution makes it a very useful distribution in practice. Generalised linear models (GLM) based on this distribution have been developed to analyse the regression effects in cross-sectional data where the count responses are collected at a single time-point<sup>[3,6]</sup>. The estimation of regression parameters in cross-sectional models can be done using maximum likelihood or quasi-likelihood estimation procedures<sup>[1,6]</sup>. Applications of GLM's are demonstrated for the analysis of over- as well as under- dispersed cross-sectional counts arising in different studies<sup>[4,7,9]</sup>. However, in various fields of research, count responses from a group of independent subjects are collected repeatedly over several time points and one is interested in estimating the effects of some known factors on these longitudinal count responses. Modeling of such longitudinal data is quite challenging as the repeated counts of each subject are correlated and the correlation pattern is unknown in practice. This is because it is extremely difficult to describe the joint distribution of correlated count observations. The problem gets further mounted up when the counts are over- or under- dispersed.

For the analysis of over-dispersed longitudinal count data, models have been developed by some authors. For example, we refer to random effects models by Thall & Vail<sup>[13]</sup> and the marginal model by Jowaheer & Sutradhar<sup>[5]</sup>. Concerning the analysis of under-dispersed longitudinal count data, there exists no literature. In this paper, we construct longitudinal Com-Poisson model which can analyse any of the three types of count responses : under-, over- or equi-dispersed. This model is basically an extension of Com-Poisson cross-sectional model developed by Jowaheer & Mamodekhan<sup>[6]</sup>. However, in

extending a model from cross-sectional to longitudinal set-up, one has to adequately model the correlations underlying the responses that are repeatedly collected over time. The longitudinal correlations can be modeled using random effects [see Thall & Vail<sup>[13]</sup>], 'working' ARMA structures [see Liang & Zeger<sup>[8]</sup>] or 'true' Gaussian type ARMA structures [see Jowaheer & Sutradhar<sup>[5]</sup>; Sutradhar<sup>[11]</sup>]. It has been shown that 'true' Gaussian type ARMA correlation structures are the most efficient and reliable in the sense that they are able to model the true underlying correlations of the repeated counts (Crowder<sup>[2]</sup>; Sutradhar & Das<sup>[12]</sup>). Developing ARMA correlation structures of different orders for correlated random variables following a discrete distribution is quite challenging and Com-Poisson distribution in particular has a complicated probability mass function. Here, we provide a method for developing and analysing Com-Poisson model with true AR(1)-type correlation structure under the longitudinal set-up.

### Com-Poisson Random Variable

Let  $y_t$  be the count response at time  $t$  ( $i = 1, 2, \dots, I; t = 1, 2, \dots, T$ ), following Com-Poisson distribution. Then

$$(1) \quad f(y_t) = \frac{\lambda^{y_t}}{(y_t!)^\nu} \frac{1}{Z(\lambda, \nu)},$$

where

$$(2) \quad Z(\lambda, \nu) = \sum_{j=0}^{\infty} \frac{\lambda^j}{(j!)^\nu}, \lambda > 0$$

and the parameter  $\nu$  is the dispersion index such that  $\nu = 1$ ,  $\nu < 1$  and  $\nu > 1$  correspond to equi-, over- and under- dispersion. Since equation (1) does not have closed form expressions, we use an asymptotic expression for  $Z(\lambda, \nu)$  proposed by Shmueli et. al<sup>[10]</sup> given by

$$(3) \quad Z(\lambda, \nu) \simeq \frac{\exp(\nu \lambda^{\frac{1}{\nu}})}{\lambda^{\frac{\nu-1}{2\nu}} (2\pi)^{\frac{\nu-1}{2}} \sqrt{\nu}}$$

The approximation is particularly good for  $\lambda > 10^\nu$ . Hence,

$$(4) \quad E(Y_t) = \theta = \lambda^{1/\nu} - \frac{\nu - 1}{2\nu}, \quad Var(Y_t) = \frac{\lambda^{1/\nu}}{\nu}$$

We can say that

$$(5) \quad Y_t \sim CMP\left(\frac{\theta}{\nu}, \nu\right)$$

As the responses are repeatedly collected over T time points,  $y_1, y_2, \dots, y_T$  are likely to be correlated. Following Sutradhar & Das<sup>[12]</sup>, these T repeated counts have Gaussian-type true but unknown auto-correlation structure given by

$$(6) \quad C_i(\rho_1, \dots, \rho_{T-1}) = \begin{pmatrix} 1 & \rho_1 & \rho_2 & \dots & \rho_{T-1} \\ \rho_1 & 1 & \rho_1 & \dots & \rho_{T-2} \\ \vdots & \vdots & \vdots & & \vdots \\ \rho_{T-1} & \rho_{T-2} & \rho_{T-3} & \dots & 1 \end{pmatrix}$$

yielding the autocovariance structure

$$(7) \quad \Sigma(\tilde{\rho}) = A^{\frac{1}{2}} C(\rho_1, \dots, \rho_{T-1}) A^{\frac{1}{2}}$$

with  $A = \text{diag } \text{var}(Y_t)$ . We now derive the components of matrix  $C(\rho_1, \dots, \rho_{T-1})$  under stationary AR(1) correlation pattern.

### Generating a Sequence of AR(1)-Correlated Com-Poisson Responses

Let response random variable  $y_{t-1}$  at  $(t-1)$ th time point follows  $CMP(\frac{\theta}{\nu}, \nu)$  and the error term  $d_t \sim CMP(\frac{(1-\rho)\theta}{\tilde{\nu}}, \tilde{\nu})$  where

$$(8) \quad \tilde{\nu} = \frac{(2q_0 + 1) + \sqrt{(2q_0 + 1)^2 - 8q_1}}{4q_1}$$

with  $q_0 = (1 - \rho)\theta$ ,  $q_1 = \frac{q_0}{\nu} [1 + \rho(1 - \nu)] + \frac{\nu-1}{2\nu^2}(1 - \rho^2)$  and  $0 < \rho < 1$ . Then, the stationary AR(1) model for the sequence  $\{y_t\}$  can be formulated as

$$(9) \quad y_t = \rho * y_{t-1} + d_t$$

where  $0 < \rho < 1$ . The symbol  $*$  indicates the binomial convolution thinning operation such that

$$(10) \quad \rho * y_{t-1} = \sum_{j=1}^{y_{t-1}} b_j(\rho) = g_t.$$

where  $\text{prob}[b_j(\rho) = 1] = \rho$  and  $\text{prob}[b_j(\rho) = 0] = 1 - \rho$ .

$$(11) \quad (g_t | y_{t-1}, \rho) \sim \text{Binomial}(y_{t-1}, \rho)$$

Hence, using equation (9), we obtain

$$(12) \quad \begin{aligned} E(Y_t) &= E_{Y_{t-1}} E(\rho * Y_{t-1} | Y_{t-1}, \rho) + E(d_t) \\ &= E(\rho Y_{t-1}) + E(d_t) \\ &= \rho E(Y_{t-1}) + (1 - \rho)\theta \\ &= \theta \end{aligned}$$

$$(13) \quad \begin{aligned} \text{Var}(Y_t) &= \text{Var}(\rho * Y_{t-1}) + \text{Var}(d_t) \\ &= \text{Var}(E(\rho * Y_{t-1})) + E(\text{Var}(\rho * Y_{t-1} | Y_{t-1})) + \text{Var}(d_t) \\ &= \text{Var}(\rho Y_{t-1}) + E(Y_{t-1} \rho(1 - \rho)) + \text{Var}(d_t) \\ &= \rho^2 \left[ \frac{\theta}{\nu} + \frac{\nu - 1}{2\nu^2} \right] + \rho(1 - \rho)\theta + \left( \frac{\theta}{\nu} + \frac{\nu - 1}{2\nu^2} \right) (1 - \rho^2) - \rho(1 - \rho)\theta \\ &= \frac{\theta}{\nu} + \frac{\nu - 1}{2\nu^2} \end{aligned}$$

$$(14) \quad \begin{aligned} E(Y_t Y_{t-k}) &= E_{Y_{t-k}} E_{Y_{t-k+1}} \dots E_{Y_t} (Y_t Y_{t-k} | Y_{t-1}, Y_{t-2}, \dots, Y_{t-k}) \\ &= E_{Y_{t-k}} [\rho^k Y_{t-k}^2 + \rho^{k-1} (1 - \rho)\theta^2 + \dots + \rho(1 - \rho)\theta^2 + (1 - \rho)\theta | Y_{t-k}] \\ &= \rho^k \left( \frac{\theta}{\nu} + \frac{\nu - 1}{2\nu^2} + \theta^2 \right) + \theta^2 (1 - \rho) (\rho^k + \rho^{k-1} + \dots + 1) \\ &= \rho^k \left[ \left( \frac{\theta}{\nu} + \frac{\nu - 1}{2\nu^2} \right) + \theta^2 \right] + \theta^2 (1 - \rho^k) \end{aligned}$$

$$(15) \quad \text{Cov}(Y_t, Y_{t-k}) = E(Y_t Y_{t-k}) - E(Y_t)E(Y_{t-k}) = \rho^k \left( \frac{\theta}{\nu} + \frac{\nu - 1}{2\nu^2} \right)$$

and

$$(16) \quad \rho_k = \text{Corr}(Y_t, Y_{t-k}) = \frac{\text{Cov}(Y_t, Y_{t-k})}{\sqrt{\text{Var}(Y_t)} \sqrt{\text{Var}(Y_{t-k})}} = \rho^k$$

Equation(16) justifies AR(1)-type stationary auto-correlations between longitudinal Com-Poisson counts modelled by equation (9).

**Longitudinal Com-Poisson Regression Model**

Let  $y_{it}$  be the count response of the  $i$ th individual at time  $t$  ( $i = 1, 2, \dots, I; t = 1, 2, \dots, T$ ). Let  $x_{it}$  be the  $p$  dimensional vector of covariates corresponding to  $y_{it}$ . Let  $\beta$  be the  $p$  dimensional vector of regression parameters. The Com-Poisson regression model can thus be formulated by replacing  $y_t$  with  $y_{it}$  and  $\lambda_t$  by  $\lambda_{it}$  in equations (1) and (2) where

$$(17) \quad \lambda_{it} = \exp(x_{it}^T \beta)$$

Hence,

$$(18) \quad f(y_{it}) = \frac{\exp(x_{it}^T \beta y_{it}) [\exp(x_{it}^T \beta (\frac{\nu-1}{2\nu})) (2\pi)^{\frac{\nu-1}{2\nu}} \sqrt{\nu}]}{(y_{it}!)^\nu [\exp(\nu \exp(\frac{x_{it}^T \beta}{\nu}))]}$$

where

$$(19) \quad E(Y_{it}) = \theta_{it} = \lambda_{it}^{1/\nu} - \frac{\nu-1}{2\nu}, \quad Var(Y_{it}) = \frac{\lambda_{it}^{1/\nu}}{\nu}$$

and the autocovariance structure,

$$(20) \quad \Sigma_i(\tilde{\rho}) = A_i^{\frac{1}{2}} C_i(\rho_1, \dots, \rho_{T-1}) A_i^{\frac{1}{2}}$$

with  $A_i = \text{diag } var(Y_{it})$ . The correlation matrix is unique for all the  $I$  clusters ,i.e,  $C_i(\rho) = C(\rho)$ .

**Generalized Quasi-likelihood Estimating Equations**

The GQL equations to estimate the regression and dispersion parameters of the longitudinal model may be written as

$$(21) \quad \sum_{i=1}^I D_i^T \widetilde{\Sigma}_i^{-1} (f_i - \mu_i) = 0,$$

where  $f_i = (f_{i1}^T, \dots, f_{it}^T, \dots, f_{iT}^T)^T$ ,  $\mu_i = (\mu_{i1}^T, \dots, \mu_{it}^T, \dots, \mu_{iT}^T)^T$  are  $2T \times 1$  vectors with  $f_{it} = (y_{it}, y_{it}^2)$ ,  $\mu_{it} = (\theta_{i1}, m_{i1})^T$ .  $\theta_{i1} = E(Y_{it})$  and  $m_{i1} = E(Y_{it}^2)$

The components  $D_i, \widetilde{\Sigma}_i, f_i$  and  $\mu_i$  are obtained in the same way as in Jowaheer & Sutradhar<sup>[5]</sup> such that

$$D_i = [\partial \mu_i / \partial \beta^T, \partial \mu_i / \partial \nu] = [D_{i1}^T, \dots, D_{it}^T, \dots, D_{iT}^T]^T,$$

with

$$D_{it} = \begin{pmatrix} \partial \theta_{i1} / \partial \beta^T & \partial \theta_{i1} / \partial \nu \\ \partial m_{i1} / \partial \beta^T & \partial m_{i1} / \partial \nu \end{pmatrix}$$

for  $t = 1, \dots, T$  where

$$\begin{aligned} \partial \theta_{i1} / \partial \beta^T &= \frac{\lambda_{i1}^{1/\nu}}{\nu} x_{i1}^T \\ \partial \theta_{i1} / \partial \nu &= \frac{\nu-1}{2\nu^2} - \frac{1}{2\nu} - \frac{\lambda_{i1}^{1/\nu} x_{i1}^T \beta}{\nu^2} \\ \partial m_{i1} / \partial \beta^T &= x_{i1}^T \left( \frac{2\lambda_{i1}^{1/\nu} + 2\nu\lambda_{i1}^{2/\nu} - \nu\lambda_{i1}^{1/\nu}}{\nu^2} \right) \\ \partial m_{i1} / \partial \nu &= \frac{1}{2\nu^3} [2\lambda_{i1}^{1/\nu} \nu \ln(\lambda_{i1}) + \nu - 1 - 4\lambda_{i1}^{2/\nu} \ln(\lambda_{i1}) \nu - 4\lambda_{i1}^{1/\nu} \nu - 4\lambda_{i1}^{1/\nu} \ln(\lambda_{i1})]. \end{aligned}$$

The covariance matrix of  $f_i$  is expressed as

$$\widetilde{\Sigma}_i = \begin{pmatrix} \widetilde{\Sigma}_{i1} & \widetilde{\Omega}_{i12} & \widetilde{\Omega}_{i13} & \cdots & \widetilde{\Omega}_{i1T} \\ & \widetilde{\Sigma}_{i2} & \widetilde{\Omega}_{i23} & \cdots & \widetilde{\Omega}_{i2T} \\ & & \widetilde{\Sigma}_{i3} & \cdots & \widetilde{\Omega}_{i3T} \\ & & & \ddots & \\ & & & & \widetilde{\Sigma}_{iT} \end{pmatrix}$$

where the diagonal submatrix

$$\widetilde{\Sigma}_{it} = \begin{pmatrix} \text{var}(Y_{it}) & \text{cov}(Y_{it}, Y_{it}^2) \\ & \text{var}(Y_{it}^2) \end{pmatrix}$$

and for  $t \neq w$ , the off-diagonal submatrix

$$\widetilde{\Omega}_{itw} = \begin{pmatrix} \text{cov}(Y_{it}, Y_{iw}) & \text{cov}(Y_{it}, Y_{iw}^2) \\ \text{cov}(Y_{it}^2, Y_{iw}) & \text{cov}(Y_{it}^2, Y_{iw}^2) \end{pmatrix}$$

for  $t = 1, \dots, T$  and  $w = 1, \dots, T$ .

The GQL estimating equation (21) is solved by the Newton-Raphson iterative method to obtain the estimates  $\hat{\beta}$  and  $\hat{\nu}$  such that

$$(22) \quad \begin{pmatrix} \hat{\beta}_{r+1} \\ \hat{\nu}_{r+1} \end{pmatrix} = \begin{pmatrix} \hat{\beta}_r \\ \hat{\nu}_r \end{pmatrix} + \left[ \sum_{i=1}^I D_i^T \widetilde{\Sigma}_i^{-1} D_i \right]_r^{-1} \left[ \sum_{i=1}^I D_i^T \widetilde{\Sigma}_i^{-1} (f_i - \mu_i) \right]_r$$

where  $\hat{\beta}_r$  is the value of  $\hat{\beta}$  at the  $r^{th}$  iteration.  $[\cdot]_r$  is the value of the expression at the  $r^{th}$  iteration.  $\rho_{|t-w|}$  is consistently estimated using the method of moments

$$(23) \quad \hat{\rho}_l = \frac{\sum_{i=1}^I \sum_{t=1}^{T-l} \tilde{y}_{it} \tilde{y}_{i,t+l} / (T-l)}{\sum_{i=1}^I \sum_{t=1}^T \tilde{y}_{it}^2 / T}$$

for  $(l = |t-w| = 1, \dots, T-1)$  where  $\tilde{y}_{it} = \frac{y_{it} - \theta_{i1}}{\sqrt{\frac{\lambda_{i1}^{\frac{1}{\nu}}}{\nu}}}$ . The algorithm works as follows: For an initial value

of  $\hat{\beta}$  and  $\hat{\nu}$ , we calculate  $\hat{\rho}_l$  using (23) and then use these two sets of parameters to update the values of  $\hat{\beta}$  and  $\hat{\nu}$ . Then the new set of parameters is used to calculate  $\hat{\rho}_l$  and the iteration continues in this way until convergence. The estimators are consistent and under mild regularity conditions, for  $I \rightarrow \infty$ , it may be shown that  $I^{\frac{1}{2}}((\hat{\beta}, \hat{\nu}) - (\beta, \nu))^T$  has an asymptotic normal distribution with mean 0 and covariance matrix  $I \left[ \sum_{i=1}^I D_i^T \widetilde{\Sigma}_i^{-1} D_i \right]^{-1} \left[ \sum_{i=1}^I D_i^T \widetilde{\Sigma}_i^{-1} (f_i - \mu_i)(f_i - \mu_i)^T \widetilde{\Sigma}_i^{-1} D_i \right] \left[ \sum_{i=1}^I D_i^T \widetilde{\Sigma}_i^{-1} D_i \right]^{-1}$

### Conclusions

In this paper, we provide a method of generating a sequence of Com-Poisson counts exhibiting Gaussian AR(1)-type autocorrelations. AR(1) autocorrelation structure has been then used to develop Com-Poisson longitudinal model to estimate the effects of fixed covariates on the repeatedly observed under-, over- and equi-dispersed count responses obtained from a large number of independent individuals. The regression and dispersion parameters are estimated by generalised quasi-likelihood estimation approach, whereas correlation parameters are estimated using method of moments.

## REFERENCES

1. Consul, P. and Shoukri, M.(1984). Maximum likelihood estimation for generalized Poisson distribution. *Communication in Statistics, Theory and Methodology*, 17, 1533-1547.
2. Crowder, M. (1995) On the use of a working correlation matrix in using generalised linear models for repeated measures. *Biometrika*, 82.
3. Guikema, S. (2007). Formulating informative, data-based priors for failure probability estimation in reliability analysis. *Reliability Engineering and System Safety*, 92(4), 490-502.
4. Guikema, S. and Gofelt, J. (2008). A flexible count data regression model for risk analysis. *Risk Analysis*, 28, 213-223.
5. Jowaheer, V. and Sutradhar, B.C. (2002). Analysing longitudinal count data with overdispersion. *Biometrika*, 89, 389-399.
6. Jowaheer, V. and Mamodekhan, N. (2009). Estimating regression effects in Com-Poisson generalized linear model. *International Journal of Mathematical and Statistical Sciences*. 1:2, 59-63.
7. Kadane, G., Shmueli, G., Minka, G., Borle, T. and Boatwright, P. (2006). Conjugate analysis of the Conway Maxwell Poisson distribution. *Bayesian Analysis*, 1, 363-374.
8. Liang, K. and Zeger, S. (1986). Longitudinal data analysis using generalised linear models. *Biometrika* 73, 13-22.
9. Lord, S., Guikema, S., Geedipally, S. (2008). Application of the Conway-Maxwell-Poisson generalised linear model for analysing motor vehicle crashes, *Accident Analysis and Prevention* 40(3), 1123-1134.
10. Shmuelli, G., Minka, T.P., Kadane, J., Borle, S., and Boatwright, P. (2005). A useful distribution for fitting discrete data: Revival of the Com-Poisson. *Applied Statistics: Journal of Royal Statistical Society, Series C*, 54(1), 127-142.
11. Sutradhar, B. C. (2003). An overview on regression models for discrete longitudinal responses. *Statistical Science*, 18(3), 377-393.
12. Sutradhar, B.C. and Das, K. (1999). On the efficiency of regression estimators in generalised linear models for longitudinal data, *Biometrika*, 86, 459-465.
13. Thall, P. and Vail, S. (1990). Some covariance models for longitudinal count data with overdispersion. *Biometrika* 46, 657-671.