

Markov Chain Monte Carlo Analysis of Underreported Count Data With an Application to Worker Absenteeism

RAINER WINKELMANN¹

Department of Economics, University of Canterbury, Christchurch, New Zealand

Abstract: A new approach for modeling under-reported Poisson counts is developed. The parameters of the model are estimated by Markov Chain Monte Carlo simulation. An application to workers absenteeism data from the German Socio-Economic Panel illustrates the fruitfulness of the approach. Worker absenteeism and the level of pay are unrelated, but absence rates increase with firm size.

JEL Classification System-Numbers: C15, C21

1 Introduction

In recent years, the available class of econometric models for analyzing count data has seen a substantial expansion (See Cameron and Trivedi (1986) and Winkelmann and Zimmermann (1995) for literature surveys). Problems like unobserved heterogeneity, selectivity, and endogeneity require, and have received, non-standard treatment when the dependent variable is a non-negative integer. An hitherto underdeveloped area of research is that of selective, or under-, reporting, a situation in which the reported number of events is less than the total number of events. A further shortcoming of previous research is the almost exclusive focus on likelihood based estimation techniques at the expense of Bayesian techniques.

With the recent advent of Markov chain Monte Carlo methods (Gelfand and Smith, 1990, Chib, Greenberg and Winkelmann, 1996), the computational problems of Bayesian methods no longer necessarily exceed those of maximum likelihood methods, and the neglect of Bayesian solution is no longer justifiable. In this paper, Markov chain Monte Carlo methods are used to provide a full Bayesian treatment of a Poisson regression model with underreporting. In addition to solving a particular problem, this paper also attempts to demonstrate that Markov chain Monte Carlo methods are simple to implement and might become a standard tool for future applied count data research.

¹ Valuable comments by John London-Lane, Siddharta Chib and two anonymous referees are gratefully acknowledged.

The interest in the Poisson model with under-reporting arose from the need to estimate the determinants of worker's absenteeism from workplace using data from the German Socio-Economic Panel. Workplace absenteeism has received increased attention recently (See Allen (1981), Barmby, Orme and Treble (1991, 1995), Delgado and Kniesner (1994), Johansson and Palme (1996)). Empirical evidence on the link between absenteeism and work place conditions and wages can provide a test of efficiency wage type explanations of unemployment (Shapiro and Stiglitz, 1984). If firms pay efficiency wages to reduce shirking, we expect the number of absent days during a given period of time to decrease as a function of pay, once other factors are accounted for. This hypothesis can be tested using micro data on counts of absent days and appropriate covariates.

The dependent variable is the number of absent days during a given period of time (the year of 1985) in a sample of 1266 blue collar workers. In contrast to previous specifications of count data models for absenteeism, I will allow for the fact that the observed number of absent days may be under-reported. There are several reasons why this might be the case. *Under-reporting* may result from an insufficient surveillance mechanism if the data are provided by the employer. Alternatively, it may be due to a lack of memory if the number of events is reconstructed retrospectively by the employee. In the dataset used in this paper, workers are typically interviewed in the middle of a calendar year and report the number of absent days during the previous year. Another possibility is that workers report spells of absenteeism only when they are linked to serious health problems. Whatever the source, under-reporting will invalidate the assumptions of the standard Poisson regression model.

I derive a modified latent Poisson regression model that addresses this situation and I use Markov chain-Monte Carlo methods to implement a data augmentation algorithm, in which the latent total number of absent days is replaced by a simulated value. With this technique point estimates and estimated standard errors for all the parameters of the model can be obtained.

2 Model

Let y_i^* denote the total number of events during a fixed time period t for individual i , and assume that y_i^* conditional on covariates x_i is Poisson distributed with mean

$$E(Y_i | x_i) \equiv \lambda_i = \exp(x_i' \beta), \quad i = 1, \dots, n \quad (1)$$

$\beta = (\beta_1, \dots, \beta_k)'$ is a vector of unknown regression coefficients and $x_i' = (x_{i1}, \dots, x_{ik})$ includes covariates such as wage, health status, and type of work. This set-up corresponds to the standard exponential Poisson regression model with

probability function

$$P(y_i^* | x_i) = \frac{\exp \{y_i^* x_i' \beta - \exp(x_i' \beta)\}}{y_i^*!} \tag{2}$$

Given a sample of n independent pairs of observations (y_i^*, x_i) , an estimate for β is readily obtained through maximum likelihood estimation or iteratively reweighted least squares as solution to

$$\sum_{i=1}^n (y_i^* - \exp(x_i' \beta)) x_i = 0 \tag{3}$$

The assumption that y^* is fully observed is unrealistic in many count data applications. Rather, it is possible that the *reported* number of events y constitutes only a fraction of all events and that count data are under-reported. Suppose that y_i , conditional on y_i^* , is binomial distributed

$$P(y_i | y_i^*, p_i) = \frac{y_i^*!}{(y_i^* - y_i)! y_i!} p_i^{y_i} (1 - p_i)^{y_i^* - y_i} \tag{4}$$

where p_i gives the individual probability of reporting an event. This probability is assumed to be constant and identical for all events and independent of the history of the process. A given number of reported events can then arise in many ways. For instance, $y_i = y_i^*$ and all the events are reported. Alternatively, $y_i = y_i^* - n$ where $n < y_i^*$ can be any number of non-reported events. The marginal distribution of the number of reported events y_i can be calculated as

$$\begin{aligned} P(Y_i = y) &= \sum_{y^* \geq y} \frac{\lambda^{y^*} e^{-\lambda}}{y^*!} \frac{y^*!}{(y^* - y)! y!} p^y (1 - p)^{y^* - y} \\ &= \sum_{y^* \geq 0} \frac{\lambda^{y^*} (1 - p)^{y^*}}{y^*!} \frac{(p\lambda)^y e^{-\lambda}}{y!} \\ &= \frac{e^{-\lambda p} (\lambda p)^y}{y!} \end{aligned} \tag{5}$$

Hence, the number of observed events is again Poisson distributed with mean $\lambda_i p_i$. An alternative derivation of this result uses the concept of probability generating functions (See, for instance, Feller 1971).

The simple form of the mixture, or compound, Poisson distribution depends on the independence assumption between the two processes. This assumption can be somewhat relaxed by letting p_i depend on a set of covariates that possibly overlap with x_i and requiring conditional independence (as in Winkelmann and Zimmermann, 1993). Here, the focus is on computational issues, and p_i is treated like a random effect.

The resulting Poisson model with underreporting is similar to a Poisson model with unobserved heterogeneity. In particular, the mean function can be rewritten as

$$E(Y_i | x_i, p_i) = \exp(x_i' \beta + \ln p_i) \tag{6}$$

Since $p_i \in [0, 1]$, the “additive error” in p_i is strictly negative. The marginal (with respect to p_i) expectation and variance of y_i are given by

$$E(Y_i|x_i) = \lambda E(p_i)$$

and

$$\begin{aligned} \text{Var}(Y_i|x_i) &= E_p(\text{Var}(Y_i|x_i, p_i)) + \text{Var}_p(E(Y_i|x_i, p_i)) \\ &= E(Y_i|x_i) + \lambda^2 \text{Var}(p_i) \end{aligned}$$

respectively. Hence, $\text{Var}(Y_i|x_i) > E(Y_i|x_i)$ and random underreporting, like unobserved heterogeneity, leads to overdispersion.

2.1 Maximum Likelihood Estimation

The p_i 's cannot be treated as parameters in a conventional sense since the resulting model is singular with $n + k$ parameters and n data points. Alternatively, the p_i 's can be treated as random effects and one could consider the marginal distribution of y_i after p_i has been integrated out. Such a procedure is common practice in models with unobserved heterogeneity. For instance, a gamma distributed multiplicative error gives rise to the negative binomial model. Here, the marginal density is given by

$$P(Y = y) = \int_0^1 \frac{e^{-\lambda p} (\lambda p)^y}{y!} f(p) dp \quad (7)$$

An algebraically convenient choice for $f(p)$ is the standard uniform distribution. The uniform leads to a closed form representation of the marginal distribution of y

$$\begin{aligned} P(Y = y) &= \int_0^1 \frac{e^{-\lambda p} (\lambda p)^y}{y!} dp \\ &= \lambda^{-1} \left(1 - e^{-\lambda} \sum_{i=0}^y \frac{\lambda^i}{i!} \right) \quad y = 0, 1, \dots \end{aligned} \quad (8)$$

In principle, more general mixing distributions could be employed based on simulated maximum likelihood techniques (Gourieroux and Monfort, 1991). Disadvantages of this approach are that the choice of a particular distribution of $f(p)$ is restrictive and difficult to justify in practice. Further, this approach is completely uninformative with regard to the p_i 's. The Bayesian approach, in contrast, yields a full set of posterior distributions for the p_i 's.

2.2 Bayesian Approach

The density of $y = (y_1, \dots, y_n)$ conditional on β and p is given by

$$P(y|\beta, p) = \prod_{i=1}^n \frac{\exp(-\exp(x'_i\beta)p_i) \exp(y_i x'_i\beta) p_i^{y_i}}{y_i!} \tag{9}$$

In Bayesian inference, we are interested in the joint posterior distribution $P(y^*, p, \beta|y, x)$ which is proportional to

$$P(y|y^*, p, \beta) \times P(y^*|\beta) \times \pi(\beta) \times \pi(p)$$

The following prior distributions π will be used:

$$\pi_\beta \sim N_k(\beta_0, B_0^{-1}) .$$

and

$$\pi_p \sim \mathcal{U}(0, 1)$$

Further, the conditional distribution of y is given in (4) and the conditional distribution of y^* in (2).

The resulting joint posterior distribution of y_i^* , p_i , and β is then proportional to

$$P(y^*, p, \beta|y, x) \propto \exp\left(-\frac{1}{2}(\beta - \beta_0)' B_0(\beta - \beta_0)\right) \times \prod_{i=1}^n \exp\{y_i^* x'_i\beta - \exp(x'_i\beta)\} \frac{p_i^{y_i}(1 - p_i)^{y_i^* - y_i}}{(y_i^* - y_i)! y_i!} \tag{10}$$

While it is intractable to analytically derive the marginal posterior distributions for the parameters of interest from (10), the MCMC approach allows to simulate the joint posterior density. Once the joint distribution is simulated, all aspects relevant to the analysis, like means and standard deviations of the marginal posteriors, can be derived.

2.3 Posterior Sampling Method

The MCMC algorithm bears close resemblance to data augmentation techniques: although y^* is latent, the conditional distribution of y^* given β is known. Hence, simulated data can replace the actual observations in an iterative sampling algorithm. Formally, the MCMC algorithm exploits the fact that, while the joint posterior (10) is complicated, the conditional distributions of a single parameter conditional on the remaining parameters are tractable. The algo-

rithm requires the successive sampling of y^* , p , and β from the following full conditional distributions:

$$y^*|p, \beta; p|y^*; \beta|y^*, p ;$$

The sampling process is initiated with values in the support of the posterior density. It is not difficult to show [by verifying the conditions of Roberts and Smith (1994)] that the Markov chain induced by this sampling process converges to the target joint posterior distribution. The sampling from the marginal distributions uses the following algorithms:

i) Sampling of p

From (10) it is seen that

$$P(p|y^*) \propto \prod_{i=1}^n p_i^{y_i^*} (1 - p_i)^{y_i^* - y_i} \tag{11}$$

$p|y^*$ is beta distributed. An algorithm for sampling from the beta distribution is given in Knuth (1969). For simulating the beta distribution it is useful to note that if X_1 and X_2 are independently gamma distributed with parameters α and ϑ , respectively, then $X_1/(X_1 + X_2)$ is beta $B(\alpha, \vartheta)$. Simulation of gamma variates, in turn, can be based on the exponential distribution since α and ϑ are integers.

*ii) Sampling of y^**

The conditional distribution of the latent count is proportional to

$$P(y^*|\beta, p) \propto \prod_{i=1}^n \frac{\exp(y_i^* x_i \beta) (1 - p_i)^{y_i^*}}{(y_i^* - y_i)!} \tag{12}$$

This is the kernel of a displaced Poisson distribution (Johnson and Kotz, 1970) where

$$P(y_i^*|\beta, p) = \frac{\{(1 - p_i) \exp(x_i \beta)\}^{y_i^* - y_i} \exp\{-(1 - p_i) \exp(x_i \beta)\}}{(y_i^* - y_i)!} \tag{13}$$

$$y_i^* = y_i, y_i + 1, \dots$$

To simulate from this distribution, Poisson random numbers are shifted by y_i .

iii) Sampling of β

Ignoring constants that do not depend on β , the full conditional density of β is

$$P(\beta|y^*) \propto \exp\left(-\frac{1}{2}(\beta - \beta_0)' B_0(\beta - \beta_0)\right) \prod_{i=1}^n \exp\{y_i^* x_i \beta - \exp(x_i \beta)\} \tag{14}$$

This density does not belong to a known family so recourse must be taken to the Metropolis-Hastings (M-H) algorithm [See Chib and Greenberg, 1995]. This may be implemented as follows. Let V_β denote the inverse of the observed information matrix in the model without underporting. Let $\beta^{(g)}$ denote the current value of β in the Markov chain. Let β^* denote a proposal value gener-

ated from the random-walk chain

$$\beta^* = \beta^{(g)} + \tau V_{\beta}^{-1/2} z$$

where τ is a specified scale factor and z is a standard normal vector. Let

$$\alpha(\beta^{(g)}, \beta^*) = \min\{\exp[\ln P(\beta^*|y, \{b_i\}) - \ln P(\beta^{(g)}|y, \{b_i\})], 1\} .$$

Then, the next item in the chain $\beta^{(g+1)}$ is obtained as follows:

$$\begin{aligned} \beta^{(g+1)} &= \beta^* \text{ with probability } \alpha(\beta^{(g)}, \beta^*) \\ &= \beta^{(g)} \text{ with probability } 1 - \alpha(\beta^{(g)}, \beta^*) \end{aligned}$$

If the rejection rate is high, τ needs to be adjusted. Experience shows that values between 0.5 and 1 achieve desirable rejection rates of 40 to 50 percent. Given that the simulations from the conditionals are set-up, the Markov Chain-Monte Carlo algorithm is implemented as follows:

Step 1. Specify starting values $y^{*(0)}$, $p^{(0)}$, $\beta^{(0)}$.

Step 2. For $g = 1, \dots, r$, simulate

$$\begin{aligned} y^{*(g+1)} &\quad \text{from } y^*|p^{(g)}, \beta^{(g)} \\ p^{(g+1)} &\quad \text{from } p|y^{*(g+1)}, \beta^{(g)} \\ \beta^{(g+1)} &\quad \text{from } \beta|y^{*(g+1)}, p^{(g+1)} \end{aligned}$$

Step 3. Set $g = g + 1$ and go to 2.

After an initial burn-in phase (of 2000 simulations, say), triplets of observations $(p^{(g)}, y^{*(g)}, \beta^{(g)})$ behave like drawings from the joint posterior (10). By making r large enough, any population characteristic, even the density itself, can be obtained to any degree of accuracy. The following empirical analysis is based on $r = 10000$ simulations. To analyze the results, the marginal posterior distributions can be plotted. Alternatively, means and standard errors of the marginals give point estimates and estimated standard errors of the parameters of interest.

3 Data and Estimation Results

The data come from the 1985 wave of the German Socio-Economic Panel. I focus on this particular year since it contains a module on subjective workplace evaluations and union membership that is of interest in the present context and not available in other years. I consider a subsample of 1266 male blue collar workers aged 25 or older. The dependent variable is the number of absent days during the twelve months period. The count variable takes values between 0

and 50 with mode 0, mean 8.0 and variance 132.2 . About 48% of all workers have taken at least one absent day.

The main issue of interest is the relationship between pay level and absenteeism. Delgado and Kniesner (1994) find for US data that pay and absenteeism are significantly and inversely related. Johansson and Palme (1994), using Swedish data, find a negative but insignificant effect of wages. The interpretation of any wage coefficient as *specific* hinges critically on the ability to sufficiently control for other factors that are correlated with both wages and absent rates. Among such factors are firm size, union status, and type of work. Firm size is important in efficiency wage explanations, since larger firms presumably have higher monitoring cost. Cost minimization will result in a lower level of monitoring and hence a lower probability that absenteeism is detected. Thus, for a given wage, absent rates should be positively related to firm size. To the extent that union member enjoy higher job security, the shirking model predicts higher absent rates for union members than for non-union members.

The nature of the work is captured by the following five dichotomous subjective evaluations. NOCAREER indicates that the job has no advancement opportunities – the cost of shirking are lower. MANUAL and HAZARD measure the nature and risks associated with the work. CONTROL and PARTICIPATION are indicators for the degree of autonomy involved in the work. Lastly, I use three indicators, ILLSOME, ILLBIG, and CHRONIC condition, to control for health status and thereby isolating proper “shirking” from absent days that are dictated by acute medical problems.

For this dataset, I estimate three models: the standard Poisson regression, the negative binomial model, and the Poisson regression model with underreporting. The Markov chain Monte Carlo run used 10000 iterations. 2000 initial observations have been deleted to allow the process to stabilize. The major advantage of the simulation approach is that it yields exact (small sample) marginal posterior distributions for all model parameters. Figures 1–4 show histograms of the sampling outcomes for the wage coefficient, one firm size effect, and the reporting probabilities for individuals 1 and 2. The posterior of the wage coefficient has mean 0.09; more than 95 percent of all simulated values lie between -0.1 and 0.3 .

The simulated reporting probabilities are concentrated between 0 and 0.1 for individual 1 with zero reported absent days, and between 0.3 and 0.5 for individual 2 with 15 reported absent days. The model does not per se imply that individuals with low counts have a low reporting probability. However, to the extent that the included regressors are insufficient to explain variations in the expected number of total absent days, unobserved heterogeneity may confound the under-reporting effect and cause a positive correlation between reported absent days and reporting probabilities. Disentangling unobserved heterogeneity and under-reporting is a task left for future research.

Next, I turn to a discussion of the regression results reported in table 1. The MCMC output is used to compute means and standard errors of the posterior distributions. Overall, the estimated coefficients are quite similar for the three

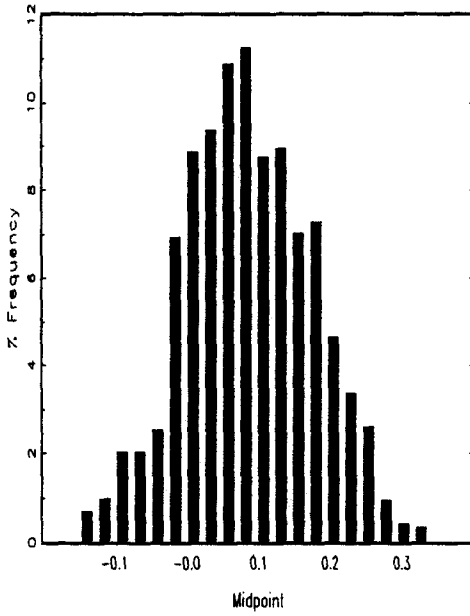


Fig. 1. Posterior for wage effect (10000 simulation)

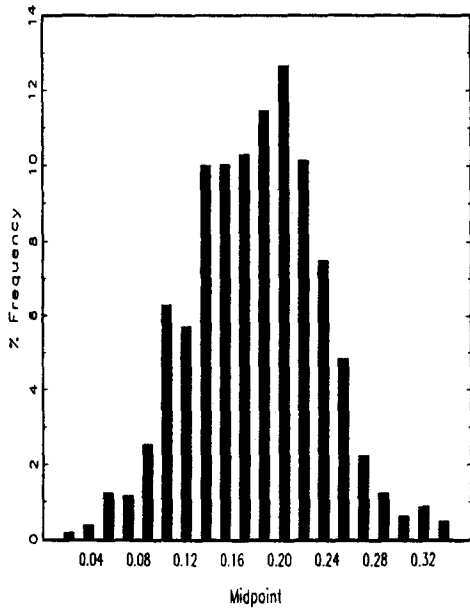


Fig. 2. Posterior for very large firm (10000 simulations)

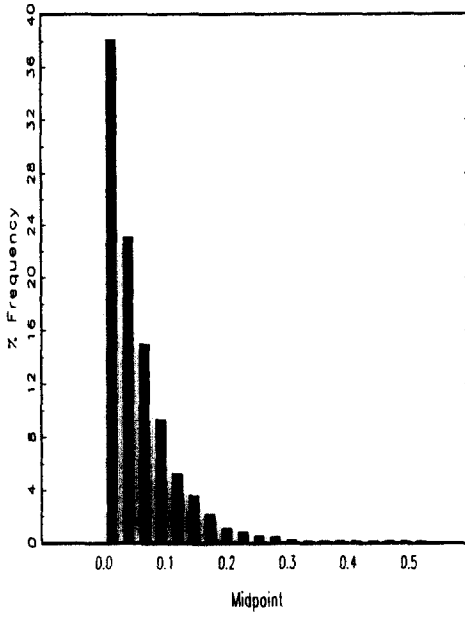


Fig. 3. Posterior for reporting probability of individual 1 ($y = 0$)

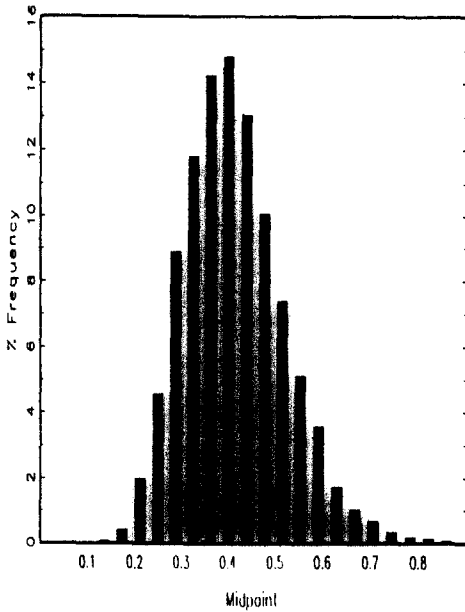


Fig. 4. Posterior for reporting probability of individual 2 ($y = 15$)

Table 1. Regression results for worker absenteeism^a (n = 1266)

	Poisson		Negative Binomial		Underreported Poisson ^b	
Constant	0.649	(0.374)	0.638	(3.643)	1.843	(0.717)
Log(wage) ^c	0.102*	(0.045)	0.091	(0.449)	0.087	(0.088)
Firm size: ^d						
21–200 empl.	0.202*	(0.033)	0.243	(0.272)	0.158*	(0.051)
201–2000 empl.	0.132*	(0.036)	0.188	(0.286)	0.180*	(0.054)
> 2000 empl.	0.475*	(0.034)	0.543*	(0.317)	0.430*	(0.058)
Health status:						
Illsome	0.446*	(0.024)	0.456*	(0.265)	0.365*	(0.050)
Illbig	0.615*	(0.037)	0.627	(0.557)	0.427*	(0.092)
Chronic Condition	0.257*	(0.025)	0.262	(0.318)	0.133*	(0.054)
Work Characteristics:						
Union	-0.112*	(0.021)	-0.111	(0.212)	-0.131*	(0.038)
No career	0.106	(0.096)	0.204	(0.962)	0.470*	(0.162)
Manual	0.070*	(0.025)	0.080	(0.240)	0.126*	(0.038)
Control	0.006	(0.021)	-0.041	(0.200)	-0.080*	(0.036)
Participation	-0.316*	(0.043)	-0.352	(0.305)	-0.165*	(0.061)
Hazard	0.037	(0.024)	0.002	(0.219)	-0.037	(0.037)
α^e			4.657	(0.275)		
Log-Likelihood	-10301.5		-3345.0			

Notes:

* Significantly different from zero at the 10 percent level.

^a Standard errors in parentheses.

^b Mean and standard error of simulated posterior distribution.

^c hourly wage.

^d Firms with 1–20 employees are reference category.

^e Overdispersion parameter. This parameter can be interpreted as the variance of a multiplicative gamma distributed error in (1). The chosen parameterization corresponds to what is known as the NEGBIN II model (See Cameron and Trivedi, 1986).

different model specifications. However, the estimated standard errors increase substantially in the under-reported Poisson model and even more in the negative binomial model. This is not surprising since it is a well known result that the Poisson standard errors are downward biased in the presence of overdispersion. But overdispersion can result both from unobserved heterogeneity (as modeled by the negative binomial model) and underreporting (as modeled by the MCMC approach).

The sign of the estimated coefficients is as expected for most variables. The number of absent days is larger for workers in larger firms. Apparently, larger firms employ lesser surveillance and, for a given wage, the expected cost of shirking are reduced. Health problems and absenteeism are positively related, whereas union members and workers with a high degree of control and participation are less absent. The same motivation effect is captured by the increased incidence of absenteeism if the job offers no career opportunities. Some of these effects are measured imprecisely and hence insignificant, in particular in the negative binomial regression.

The interesting substantive result is that the point estimate for the wage effect is positive throughout, and insignificant in the negative binomial and the under-reported Poisson models. It is safe to conclude that in this data on a sample of German blue collar workers, there is no evidence for a negative relation between pay and absenteeism. This finding complements previous results based on a different sample by Johansson and Palme (1996).

4 Conclusion

The purpose of this article is twofold. Firstly, a new count data model for underreported counts is developed. Secondly, Markov Chain Monte Carlo techniques are used to estimate posterior densities and parameters. The successful application of this recent technique in this particular model suggests that it might prove a useful tool in other count data problems as well. In an application to worker absenteeism data from the German Socio-Economic Panel I find that worker absenteeism and the level of pay are unrelated, while the number of absent days increases with firm size. A natural extension of the model to be pursued in future research will formulate and estimate a latent negative binomial model that allows for unobserved heterogeneity in addition to under-reporting.

References

- Albert J (1992) A Bayesian analysis of a Poisson random effects model for home run hitters. *The American Statistician* 46:246–253
- Allen SG (1981) An empirical model of work attendance. *Review of Economics and Statistics* 63:77–82
- Barnby TA, Orme CD, Treble JD (1991) Worker absenteeism: An analysis using micro data. *Economic Journal* 101:214–229
- Barnby TA, Orme CD, Treble JD (1995) Worker absence histories: A panel data study. *Labour Economics* 2:53–66
- Cameron AC, Trivedi PK (1986) Econometric models based on count data: Comparisons and applications of some estimators and tests. *Journal of Applied Econometrics* 1:29–53
- Chib S, Greenberg E (1995) Understanding the Metropolis-Hastings algorithm. *American Statistician* 49:327–335
- Chib S, Greenberg E, Winkelmann R (1996) Posterior simulation and model choice in longitudinal generalized linear models. Discussion Paper No. 9605, Department of Economics, University of Canterbury
- Delgado MA, Kniesner TJ (1994) Count data models with variance of unknown form: An application to a hedonic model of worker absenteeism. Universidad Carlos III de Madrid Working Paper No. 94-49

- Feller W (1971) An introduction to probability theory and its applications Vol. 2. 2nd ed., John Wiley, New York
- Gelfand AE, Smith AFM (1990) Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association* 85:398–409
- Geweke J (1995) Monte Carlo simulation and numerical integration. forthcoming in: Amman H, Kendrick D, Rust J (eds.) *Handbook of computational economics*, North Holland
- Gourieroux C, Monfort A (1991) Simulation based inference in models with heterogeneity. *Annales d'Economie et de Statistique* 20/21:69–107
- Johansson P, Palme M (1996) Do economic incentives affect work absence? Empirical evidence using Swedish micro data. *Journal of Public Economics* 59:195–218
- Johnson NL, Kotz S (1970) *Continuous univariate distributions* Vol. 1. John Wiley, New York
- Knuth DE (1969) *The art of computer programming* Vol. 2: Seminumerical algorithms. Addison Wesley Reading, Mass
- Roberts GO, Smith AFM (1992) Some convergence theory for Markov chain Monte Carlo. manuscript
- Shapiro C, Stiglitz JE (1984) Equilibrium unemployment as a worker discipline device. *American Economic Review* 74:433–444
- Winkelmann R, Zimmermann KF (1993) Poisson-logistic regression. Discussion Paper No. 93-18, Department of Economics, University of Munich
- Winkelmann R, Zimmermann KF (1995) Recent developments in count data modeling: Theory and applications. *Journal of Economic Surveys* 9:1–24

First version received: September 1995

Final version received: February 1996