

Research Article

Econometric Modelling Based on Dynamic Count Regression and China Power Supply Dataset

Yixin Yan ¹, Jiliang Hu ¹, Xiding Chen ², and A. P. Senthil Kumar ³

¹School of Economics and Business Administration, Central China Normal University, Wuhan, China

²School of Finance and Trade, Wenzhou Business College, Wenzhou, China

³School of Social Work, Jigjiga University, Somali Regional State, Jigjiga, Ethiopia

Correspondence should be addressed to Jiliang Hu; jilianghucentral@outlook.com and A. P. Senthil Kumar; senthilapsk@gmail.com

Received 18 March 2022; Revised 13 April 2022; Accepted 4 May 2022; Published 28 May 2022

Academic Editor: Amandeep Kaur

Copyright © 2022 Yixin Yan et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Traditionally, economic data of power supply is often analyzed through the count regression model due to the type of empirical data in the decision-making process. However, in reality, it is difficult to use count data model for data with autoregressive features. The main reason is that the time series features and autoregressive attributes cannot be controlled through the count regression model, which violates the assumptions set by the model. Therefore, there may be errors in the empirical analysis results. This letter firstly describes the characteristic of the count regression model and the problem, and then we refine the multiplicative autoregressive count model for dynamic count data. The model has desirable theoretical properties and is trivial to incorporate into existing models for the count data. In this study, the multiplicative autoregressive counting model for dynamic counting data is improved. The model has ideal theoretical properties and can be easily incorporated into existing economic models of counting data, especially for power supply policy analyses.

1. Introduction

In traditional economic modelling in energy problems, the traditional regression models are not suitable for analyzing nonnegative count data from an empirical dataset [1]. This is because traditional regression uses the ordinary least squares (OLS), which assumes that data are normally distributed, where count data do not follow a normal distribution [2]. Therefore, in terms of analyzing count data, three types of count regression models are often used, that is, Poisson Regression [3], Zero-Inflated Poisson Regression [4], and Negative Binominal Regression [5]. For example, Bai et al. [6] employed data from scientific research projects in China, where the research object is the number of thermodynamic research projects. Another example is that in the current COVID-19 outbreak, the number of infections among different people diagnosed is actually similar to a pattern of count data, but this pattern is more similar to a pattern with time series characteristics. The data are thus nonnegative statistical values, and it is more suitable to apply Poisson

Regression for analysis. In recent years, with the progress of measurement methods and software technology, great progress has been made in measurement problems and carding model construction that could not be solved in the past [7–9].

In energy planning, power supply policy plays a vital role in the national economy and is important for people's livelihoods. Reference [10] pointed out that during the annual peak electricity consumption period, the coordination of energy supply in China becomes a challenge. The power consumption is directly determined by the size and load of the power generation capacity. Therefore, it is important for the power sector in China to accurately predict the development of power generation in China so that policies can be made accordingly based on the data. However, as the authors in [11] pointed out, power generation in China has shown the characteristics of being long term, seasonal, and periodical. In terms of the developing trend, as time goes by and with the improvement of people's living standards, the power generation fluctuates, with an

overall trend of increase in the total amount. As for periodicity, due to the influence of various factors such as seasons or demand, the power generation in China shows a repeated cyclic change every 12 months. Reference [12] found that this phenomenon led to a serious problem that the power generation of the power plant in each period actually has a certain degree of time series problem, which also results in the existence of time series in the data used for power supply economic planning and policy formulation.

The stability of the time series directly affects the validity of the constructed model [13]. The numerical characteristics of a stable time series, such as the mean, variance, and covariance, do not change with time, and the randomness of the time series at each time point follows a certain probability distribution [14]. That is, it can use the information from a previous time point in the time series to build a model to fit the past information and then make predictions [15]. However, for a nonstable time series, the numerical characteristics change with time, and thus, it is impossible to grasp the overall randomness of the time series through known information [13]. Hence, it is difficult to model and predict a nonstable series. When a regression model is built with two independent nonstable time series, a statistically significant regression function is more likely to occur, which is called a spurious regression. The commonly used method for testing the stability of time series is the unit root test [16]. However, the unit root test is often misused with an inappropriate method, leading to wrong tests and spurious regression. Therefore, when analyzing data with an econometric model, a specific model must be selected according to the characteristics of data [17, 18]. In the face of count data with nonnegative values, it is recommended to use the count regression model. However, in reality, the observation data sometimes have the nature of time series or autoregression. In this case, analyzing data with the count regression model will cause biases in estimation because it cannot deal with the problem of autoregression and possible heteroscedasticity. When it comes to data analysis in electric power research, data in this field often have a certain degree of timing problem in the power generation of power plants in each period. If this phenomenon cannot be effectively solved, other robust analyses may still be required for empirical analysis to ensure the correctness of the empirical results. Scholars have used generalized linear models and dynamic regression models to deal with the above problems. However, using a combination of the two has proven difficult, especially because of the difficulties in using lagged dependent variables in generalized linear models as means of modelling autoregressive dynamic processes. These problems exist in dynamic count models, including the problem that taking lagged dependent variables as a regression variable often leads to econometric difficulties [19]. For this reason, the literature has proposed different approximations, including some original approaches [20, 21]. But most of these models are difficult to implement and have unrealistic assumptions about the data-generating process of the dependent variable. To address this issue, we used the improved multiplicative autoregressive model for count data proposed by [22] and proposed a revised model for the count

regression model. This paper is structured as follows: Section 2 constructs the count regression model used in traditional econometrics and explains related issues. In Section 3, a revised model incorporating the time series features and autoregressive properties is proposed. In Section 4, the paper concludes with a summary of the main findings and points out the potential application of the revised model.

2. Constructing Count Regression Models

Some statistical data follow a distribution skewed to the right and thus do not appear to be normally distributed. In this case, the classical regression model is usually not an optimal option for model estimation. In event analysis, for example, if the observed variable has a large number of values such as 0 or 1 in the data, Poisson regression is generally used for analysis [23]. This model is set as follows.

Suppose that the dependent variable representing the number of events Y_i , $i = 1, 2, \dots, n$ is a nonnegative random variable, with the observation y_i coming from the Poisson distribution $\text{Poi}(u_i)$ of the parameter λ_i , and is affected by other independent variables $\mathbf{x}_i = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1})$. The probability density function is as follows:

$$\Pr(Y_i = y_i | x_i) = \frac{\lambda_i^{y_i} \exp(-\lambda_i)}{y_i!}, \quad y_i = 0, 1, 2, \dots, u_i > 0. \quad (1)$$

In this equation, $\lambda_i > 0$ means that the Poisson arrival rate is affected by the independent variable \mathbf{x}_i . Assuming that the expected value is equal to the variance; that is, $E(Y_i = y_i | \mathbf{x}_i) = \text{Var}(Y_i = y_i | \mathbf{x}_i)$. In the Poisson distribution, λ represents the average number of events occurring per unit time. As such, the Poisson regression can be defined as follows:

$$\ln[E(Y_i = y_i | \mathbf{x}_i)] = \ln(\lambda_i) = \mathbf{x}_i' \boldsymbol{\beta}, \quad i = 1, 2, \dots, n. \quad (2)$$

Assuming that the samples are mutually independent, the likelihood function of the samples can be developed as follows:

$$L(\boldsymbol{\beta}) = \frac{\exp(-\sum_{i=1}^n \lambda_i) \prod_{i=1}^n \lambda_i^{y_i}}{\prod_{i=1}^n y_i!}. \quad (3)$$

Take the logarithm of equation (3) and we can obtain the following equation:

$$\begin{aligned} \ln L(\boldsymbol{\beta}) &= \sum_{i=1}^n [-\lambda_i + y_i \ln \lambda_i - \ln(y_i!)] \\ &= \sum_{i=1}^n [-\exp(\mathbf{x}_i' \boldsymbol{\beta}) + y_i \mathbf{x}_i' \boldsymbol{\beta} - \ln(y_i!)]. \end{aligned} \quad (4)$$

Perform first-order differentiation on equation (4) and we obtain the following:

$$\sum_{i=1}^n [y_i - \exp(\mathbf{x}_i' \boldsymbol{\beta})] \mathbf{x}_i = 0. \quad (5)$$

From the above equation, the estimated coefficient $\boldsymbol{\beta}$ of each independent variable \mathbf{x}_i with consistency can be deduced. In the above equation, X is an independent variable.

In analyzing empirical data, the influence on the dependent variable is discussed by including various micro- and macrovariables.

In addition, the limitation of Poisson regression is that the mean (expected value) and variance of the distribution are assumed to be equal, resulting in an equidispersion. However, in reality, empirical data often do not have expected values that are equal to variances. To solve this problem, equation (2) can be rewritten as follows:

$$\ln(\lambda_i) = \mathbf{x}_i \boldsymbol{\beta} + \varepsilon_i, \mathbf{i} = 1, 2, \dots, \mathbf{n}. \quad (6)$$

In this equation, ε_i represents the unobservable part or the heterogeneity of the individual in the conditional expectation, and the following equation can be obtained:

$$\lambda_i = \exp(\mathbf{x}_i \boldsymbol{\beta}) \exp(\varepsilon_i) = \mathbf{u}_i \mathbf{v}_i. \quad (7)$$

\mathbf{u}_i is the deterministic function of the independent variable, and \mathbf{v}_i is still a random variable, which belongs to the unobservable part. Therefore, the empirical model still follows the Poisson distribution:

$$P(\mathbf{Y}_i = y_i | \mathbf{x}_i, \mathbf{v}_i) = \frac{\exp(-\mathbf{u}_i \mathbf{v}_i) (\mathbf{u}_i \mathbf{v}_i)^{y_i}}{y_i!}. \quad (8)$$

As \mathbf{v}_i is usually larger than 0, it is assumed that it follows the Gamma distribution, and it is assumed that $\mathbf{v}_i \sim \text{Gamma}(1/a, a)$ and $a > 0$. Therefore, the expected value of \mathbf{v}_i is 1 and the variance is a . Substitute its probability density into the equation, a negative binomial distribution can be obtained. Perform MLB estimation, and we can get a negative binomial regression. The conditional variance is as follows:

$$\text{Var}(Y_i | x_i) = u_i + \alpha u_i^2 > u_i = E(Y_i | x_i). \quad (9)$$

This shows that the conditional variance is larger than expected in the negative binomial regression. The above equation shows that the variance value is an increasing function of α . When α approaches 0, the negative binomial regression and Poisson regression tend to be equal. In addition, if α in the Gamma distribution is transformed into δ/u_i ; that is, $\mathbf{v}_i \sim \text{Gamma}(u_i/\delta, \delta/u_i)$, then equation (9) can be turned into the following equation:

$$\text{Var}(Y_i | x_i) = u_i + \left(\frac{\delta}{u_i}\right) u_i^2 = u_i + \delta u_i > u_i. \quad (10)$$

Through the estimation of the approximate function, the analysis can be carried out according to the properties of the empirical data.

If the empirical data contains a large number of 0 values, traditionally the Zero-Inflated Poisson Regression or the Zero-Inflated Negative Binomial Regression is often used. Therefore, the empirical data also assumed that followed the Poisson distribution or negative binomial distribution. Theoretically, the analysis is divided into two parts. First, the decision to take the proportion of 0 or 0 is equivalent to binary choice in econometrics. Second, decide which value group to choose. Assume that the dependent variables follow a mixed distribution:

$$\begin{cases} P(y_i = 0 | x_i) = \theta, \\ P(y_i = j | x_i) = \frac{(1 - \theta) \exp(-\lambda_i) \lambda_i^j}{j! (1 - \exp(-\lambda_i))}, \end{cases} (j = 1, 2, \dots), \quad (11)$$

where $\lambda_i = \exp(\mathbf{x}_i \boldsymbol{\beta})$. Similarly, the estimated value of the coefficient of each independent variable x can be further obtained through calculation, and whether it has statistical significance for subsequent analysis can be decided.

Traditionally, when is it appropriate to use the Poisson regression? When is it appropriate to use negative binomial regression? Generally, in addition to looking at the data attribute, it is also a recommended practice to use respective regression analyses and calculate robust standard errors to achieve consistent estimation. This is somewhat similar to the regression estimation using the least-squares method and then getting robust standard error. Generally speaking, these approaches mostly aim to ensure that the estimation results are robust and reliable.

3. Dynamic Count Regression Models

To summarize, this paper first reviewed what kind of theoretical model and approaches are usually used when encountering count data in the process of data analysis. However, these models may produce biased results in estimation when used for analyzing data with time series or autoregressive characteristics. For example, in analyzing the problem of heat energy consumption generated in the production activities of various industrial plants, the heat energy consumption during each working day is related to the production activity plan [24]. As a result, there is a time series phenomenon according to the production planning activities. Failure to take this phenomenon into account in the analysis could lead to biased estimation. Admittedly, when using power supply data of plants, regional differences and autocorrelation issues seem to be negligible in empirical analysis if the research area is Japan, Taiwan, and other small island economies. However, when researchers aim to investigate economics with a large land area, such as China and the United States, it is important to know that these economics have different climatic zones and geographical regions. Therefore, in studying the power supply problem of the power plant, ignoring the impact of time series on the region in data analysis would lead to increased risks in model misspecification and biased estimation. Therefore, some suggestions are provided for researchers. First, the research objects must be limited, such as focusing only on the coastal areas of China or cities on the east coast of the United States, to avoid the potential impact of large variances within a large study area. Second, for data with time series features, an adjustment in the model is required. In order to model the autoregressive process while avoiding the above problems, the model can be rewritten in a way that the logarithm of the lagged dependent variable is taken as the dependent variable. In the regression, there is an infinite number of independent variables and error terms with geometric decay lags. These lag terms are in

exponential functions, constraining the conditional mean to reasonable nonnegative values. Therefore, equation (5) can be rewritten as follows:

$$\begin{aligned}
 y_t &= \exp(a_0 + b_0 x_t + \mathbf{u}_t) y_{t-1}^{a_1}, \\
 y_t &= \exp(a_0 + b_0 x_t + \mathbf{u}_t + a_1 \log(y_{t-1})), \\
 y_t &= \exp(a_0 + b_0 x_t + \mathbf{u}_t + a_1 \log(\exp(a_0 + b_0 x_{t-1} + \mathbf{u}_{t-1})) \\
 &\quad + a_1 \log(y_{t-2})), \\
 y_t &= \exp(a_0 + b_0 x_t + \mathbf{u}_t + a_0 a_1 + a_1 u_{t-1} + a_1^2 a_0 + a_1^2 b_0 x_{t-2} \\
 &\quad + a_1^2 u_{t-2} + \dots).
 \end{aligned} \tag{12}$$

This generalization of autoregression has a number of desirable properties. The model is stable for $|a_1| < 1$. After simplification, we can use the long-run multiplier for an independent variable $LRM = (b_0 + b_1 + b_2 + \dots / 1 - a_1 - a_2 - \dots)$. Although this model is desirable, it encounters problems with data where the observation count is 0. Since the logarithm of 0 cannot be calculated, this value will result in model nonconvergence in model estimation. Theoretically, the data generation process in equation (12) assumes that zero is a convergent state, once a value of zero is observed in a series, the series will continue to remain zero. However, this is not the case in all settings. Therefore, this model has to be abandoned in favour of a more flexible model. Applying the concept of Zeger and Qaqish (1998), we add a constant to the time series so that no value is equal to 0 ($y_{ts}^+ = y_{ts} + c$). We also add a constant to the observed zero values of the time series ($y_{ts}^+ = \max(y_{ts}, c, c \leq 1$). Then equation (12) can be described as follows:

$$y_t = \exp(a_0 + b_0 x_t + \mathbf{u}_t + a_1 \log(y_{t-1}^+)). \tag{13}$$

Compared with equation (12), equation (13) produces an econometric model that can still proceed to model estimation when the data contain values of 0. In this case, the count regression can be performed on the data with autocorrelation characteristics according to equation (13). Secondly, to further analyze the autoregressive problem, equation (12) can be further rewritten as follows:

$$y_t = \exp(a_0 + b_0 x_t + \mathbf{u}_t + a_1 \log(y_{t-1}^+) + \theta d_{t-1}), \tag{14}$$

where $d_{t-1} = 1$ if $y_{t-1} = 0$, $d_{t-1} = 0$ if otherwise. Equation (14) takes into account the intertemporal influence on autoregression, so the count regression model can be further extended to statistical data analysis with time series characteristics and autoregressive properties with equations (13) and (14).

4. Conclusions

In reality, observed data often have the nature of time series or autoregression. In this case, analyzing with count regression models will result in biased estimation, because the count regression models cannot handle the problem of autoregression and possible heteroscedasticity. Therefore, this study proposed an improved autoregressive model for

dynamic count data. This model has advantages over the Poisson Regression, Zero-Inflated Poisson Regression, and Negative Binominal Regression used in general econometric analysis, which cannot take into account time series and autoregression. The regression model proposed in this paper points out that it is possible to apply a link function to a lagging dependent variable and use it in the generalized linear model, which further expands this model to analyze various types of data. In this study, the revised econometric model has desirable theoretical properties and can be easily incorporated into existing models used for analyzing count data, especially for some specific policy analyses.

At the end of the article, we explain the application of the revised model. As mentioned above, to analyze the heat energy consumption generated in the production activities of each industrial plant, it is important to consider that the heat energy consumption during each working day is related to the production activity plan, so that the production planning activities have time series. Through the revised econometric model, we can further discuss how to model thermal energy data with the time series nature and analyze the influence of independent variables on the dependent variables. Secondly, the model setting of this paper can be further combined with the panel data model in the future to discuss whether the single root problem should be dealt with in data under the control of regional difference attributes and random variable characteristics under the dynamic situation to facilitate empirical application or model construction. In addition, in analyzing policies and research projects, such as thermal policies, the revised model can be extended to discuss the impact of changes in the number of projects in different regions and different periods. For instance, it can be used to analyze the thermal energy dispersal status of incinerators and power plants. Furthermore, the influence of the marginal rate of change of each time on the thermal energy dissipated scale can be analyzed under the first-order difference. In terms of empirical analysis, the revised model can also further discuss power supply issues in large regions such as China, Russia, and the United States and can be used to conduct follow-up analysis.

Data Availability

Requests for access to the data used to support the findings of this study should be made to Jiliang Hu (jilianghucentral@outlook.com).

Conflicts of Interest

The authors declare that they have no conflicts of interest regarding the publication of this paper.

Authors' Contributions

Yixin Yan and Jiliang Hu performed conceptualization, developed the methodology, and wrote the original draft. Xiding Chen performed the formal analysis. A.P.Senthil Kumar performed the formal analysis, wrote the final draft, and performed the supervision.

References

- [1] G. K. F. Tso and K. K. W. Yau, "Predicting electricity energy consumption: a comparison of regression analysis, decision tree and neural networks," *Energy*, vol. 32, no. 9, pp. 1761–1768, 2007.
- [2] E. W. Steyerberg, A. J. Vickers, N. R. Cook et al., "Assessing the performance of prediction models," *Epidemiology*, vol. 21, no. 1, pp. 128–138, 2010.
- [3] R. A. Sian and C.-C. Wang, "A generalized log-linear Poisson-modeled correlation to predict the optimal heat rejection pressure of transcritical CO₂ systems," *Science and Technology for the Built Environment*, vol. 24, no. 8, pp. 897–907, 2018.
- [4] Y. Zhang, J. Dai, B. Chen, and K. Chen, "Correlation between economic and industrial demand and scientific innovation: a case study of thermodynamics discipline statistics of National Natural Science Foundation of China," *Journal of Thermal Analysis and Calorimetry*, vol. 144, no. 6, pp. 2347–2355, 2021.
- [5] S. Li and Q. Shao, "Exploring the determinants of renewable energy innovation considering the institutional factors: a negative binomial analysis," *Technology in Society*, vol. 67, Article ID 101680, 2021.
- [6] Y. Bai, L. Chou, and W. Zhang, "Industrial innovation characteristics and spatial differentiation of smart grid technology in China based on patent mining," *Journal of Energy Storage*, vol. 43, Article ID 103289, 2021.
- [7] F. Castellares, S. L. P. Ferrari, and A. J. Lemonte, "On the Bell distribution and its associated regression model for count data," *Applied Mathematical Modelling*, vol. 56, pp. 172–185, 2018.
- [8] E. Altun, "A new model for over-dispersed count data: Poisson quasi-Lindley regression model," *Mathematical Sciences*, vol. 13, no. 3, pp. 241–247, 2019.
- [9] M. El-Morshedy, E. Altun, and M. S. Eliwa, "A new statistical approach to model the counts of novel coronavirus cases," *Mathematical Sciences*, vol. 16, no. 1, pp. 37–50, 2021.
- [10] G. Luo, X. Zhang, S. Liu, E. Dan, and Y. Guo, "Demand for flexibility improvement of thermal power units and accommodation of wind power under the situation of high-proportion renewable integration-taking North Hebei as an example," *Environmental Science and Pollution Research*, vol. 26, no. 7, pp. 7033–7047, 2019.
- [11] S. Han, Y.-h. Qiao, J. Yan, Y.-q. Liu, L. Li, and Z. Wang, "Mid-to-long term wind and photovoltaic power generation prediction based on copula function and long short term memory network," *Applied Energy*, vol. 239, pp. 181–191, 2019.
- [12] D. Yang and Z. Dong, "Operational photovoltaics power forecasting using seasonal time series ensemble," *Solar Energy*, vol. 166, pp. 529–541, 2018.
- [13] Y. Zou, R. V. Donner, N. Marwan, J. F. Donges, and J. Kurths, "Complex network approaches to nonlinear time series analysis," *Physics Reports*, vol. 787, pp. 1–97, 2019.
- [14] C.-L. Liu, W.-H. Hsiao, and Y.-C. Tu, "Time series classification with multivariate convolutional neural network," *IEEE Transactions on Industrial Electronics*, vol. 66, no. 6, pp. 4788–4797, 2019.
- [15] G. L. Simpson, "Modelling palaeoecological time series using generalised additive models," *Frontiers in Ecology and Evolution*, vol. 6, p. 149, 2018.
- [16] W. Fan and Y. Hao, "An empirical research on the relationship amongst renewable energy consumption, economic growth and foreign direct investment in China," *Renewable Energy*, vol. 146, pp. 598–609, 2020.
- [17] M. B. Shrestha and G. R. Bhatta, "Selecting appropriate methodological framework for time series data analysis," *The Journal of Finance and Data Science*, vol. 4, no. 2, pp. 71–89, 2018.
- [18] B. K. Sovacool, J. Axsen, and S. Sorrell, "Promoting novelty, rigor, and style in energy social science: towards codes of practice for appropriate methods and research design," *Energy Research & Social Science*, vol. 45, pp. 12–42, 2018.
- [19] P. T. Brandt, J. T. Williams, B. O. Fordham, and B. Pollins, "Dynamic modeling for persistent event-count time series," *American Journal of Political Science*, vol. 44, no. 4, pp. 823–843, 2000.
- [20] P. T. Brandt and J. T. Williams, "A linear Poisson autoregressive model: the Poisson AR(p) model," *Political Analysis*, vol. 9, no. 2, pp. 164–184, 2001.
- [21] T. L. Davis, B. Dirks, E. A. Carnero et al., "Chemical oxygen demand can be converted to gross energy for food items using a linear regression model," *Journal of Nutrition*, vol. 151, no. 2, pp. 445–453, 2021.
- [22] A. C. Cameron and P. K. Trivedi, *A Companion to Theoretical Econometrics*, p. 331, Wiley, New Jersey, USA, 2001.
- [23] W. Chen, L. Qian, J. Shi, and M. Franklin, "Comparing performance between log-binomial and robust Poisson regression models for estimating risk ratios under model misspecification," *BMC Medical Research Methodology*, vol. 18, no. 1, pp. 63–12, 2018.
- [24] J. S. Randhawa and I. S. Ahuja, "An investigation into manufacturing performance achievements accrued by Indian manufacturing organization through strategic 5S practices," *International Journal of Productivity and Performance Management*, vol. 67, no. 4, pp. 754–787, 2018.