MDPI

*Article*

# The Naive Estimator of a Poisson Regression Model with a Measurement Error

## Kentarou Wada * and Takeshi Kurosawa

Department of Applied Mathematics, Tokyo University of Science, Kagurazaka 1-3, Shinjuku-ku,
Tokyo 1628601, Japan
*   Correspondence: wadaken5269@gmail.com

**Abstract:** We generalize the naive estimator of a Poisson regression model with a measurement error as discussed in Kukush et al. in 2004. The explanatory variable is not always normally distributed as they assume. In this study, we assume that the explanatory variable and measurement error are not limited to a normal distribution. We clarify the requirements for the existence of the naive estimator and derive its asymptotic bias and asymptotic mean squared error (MSE). The requirements for the existence of the naive estimator can be expressed using an implicit function, which the requirements can be deduced by the characteristic of the Poisson regression models. In addition, using the implicit function obtained from the system of equations of the Poisson regression models, we propose a consistent estimator of the true parameter by correcting the bias of the naive estimator. As illustrative examples, we present simulation studies that compare the performance of the naive estimator and new estimator for a Gamma explanatory variable with a normal error or a Gamma error.

**Keywords:** Poisson regression model; error in variable; naive estimator; asymptotic bias

## 1. Introduction

We often cannot measure explanatory variables correctly in regression models because an observation may not be performed properly. The estimation result may be distorted when we estimate the model from data with measurement errors. We call models with measurement errors in an explanatory variable Error in Variable (EIV) models. In addition, actual phenomena often cannot be explained adequately by a simple linear structure, and the estimation of non-linear models, especially generalized linear models, from data with errors is a significant problem. Various studies have focused on non-linear EIV models (see, for example, Box 1963; Geary 1953). Classical error models assume that an explanatory variable is measured with independent stochastic errors (Kukush and Schneeweiss 2000). Berkson error models assume that the explanatory variable is a controlled variable with an error and that only the controlled variable can be measured (Burr 1988; Huwang and Huang 2000). Approaches to EIV models vary according to the situation. In this paper, we consider the former EIV. The corrected score function in Nakamura (1990) has been used to estimate generalized linear models. In particular, the Poisson regression model is easy to handle analytically in generalized linear models as we see later. Thus, we focus on the Poisson regression model with measurement errors.

Approaches to a Poisson regression model with classical errors have been discussed by Kukush et al. (2004), Shklyar and Schneeweiss (2005), Jiang and Ma (2020), Guo and Li (2002), and so on. Kukush et al. (2004) described the statistical properties of the naive estimator, corrected score estimator, and structural quasi score estimator of a Poisson regression model with normally distributed explanatory variable and measurement errors. Shklyar and Schneeweiss (2005) assumed an explanatory variable and a measurement error with a multivariate normal distribution and compared the asymptotic covariance matrices of the corrected score estimator, simple structural estimator, and structural quasi score estimator of a Poisson regression model. Jiang and Ma (2020) assumed a high-dimensional

explanatory variable with a multivariate normal error and proposed a new estimator for a Poisson regression model by combining Lasso regression and the corrected score function. Guo and Li (2002) assumed a Poisson regression model with classical errors and proposed an estimator that is a generalization of the corrected score function discussed in Nakamura (1990) for generally distributed errors; they derived the asymptotic normality of the proposed estimator.

In this study, we generalize the naive estimator discussed in Kukush et al. (2004). They reported the bias of the naive estimator, however, the explanatory variable is not always normally distributed as they assume. In practice, the assumption of a normal distribution is not realistic. Here, we assume that the explanatory variable and measurement error are not limited to normal distributions. However, the naive estimator does not always exist in every situation. Therefore, we clarify the requirements for the existence of the naive estimator and derive its asymptotic bias. The constant vector to which the naive estimator converges in probability does not coincide with the unknown parameter in the model. Therefore, we propose a consistent estimator of the unknown parameter using the naive estimator. It is obtained from a system of equations that represent the relationship between the unknown parameter and constant vector. As illustrative examples, we present explicit representations of the new estimator for a Gamma explanatory variable with a normal error or a Gamma error.

In Section 2, we present the Poisson regression model with measurement errors and the definition of the naive estimator and show that the naive estimator has an asymptotic bias for the true parameter. In Section 3, we consider the requirements for the existence of the naive estimator and derive its asymptotic bias and asymptotic mean squared error (MSE) assuming that the explanatory variable and measurement error are generally distributed. In addition, we introduce application examples of a Gamma explanatory variable with a normal error or a Gamma error. In Section 4, we propose the corrected naive estimator as a consistent estimator of the true parameter under general distributions and give application examples for a Gamma explanatory variable with a normal error or a Gamma error. In Section 5, we present simulation studies that compare the performance of the naive estimator and corrected naive estimator. In Section 6, we apply the naive and corrected naive estimators to real data in two cases. Finally, discussions are presented in Section 7.

## 2. Preliminary

In this section, we state the statistical model considered in this paper and the definition of the naive estimator and show that the naive estimator has an asymptotic bias for the true parameter.

### 2.1. Poisson Regression Models with an Error

We assume a single covariate Poisson regression model between the objective variable $Y$ and explanatory variable $X$

$$Y|X \sim Po(\exp(\beta_0 + \beta_1 X)).$$

$X$ can typically be correctly observed. We assume here that $X$ has a stochastic error $U$ as

$$W = X + U,$$

where $U$ is supposed to be independent of $(X, Y|X)$. We also assume that

$$(Y_i, X_i, U_i) \ (i = 1, \dots, n) \tag{1}$$

are independent and identically distributed samples of the distributions of $(Y|X, X, U)$. Although we can observe $Y|X$ and $W$, we assume that $X$ and $U$ cannot be directly observed. However, even if we know the family of the distributions of $X$ and $U$, we can-not make a

statistical inference regarding $X$ and $U$ if we can observe only $W$. Because $U$ is the error distribution, the mean of $U$ is often zero, and we may suppose that we have empirical information about the degree of error (the variance of $U$). Therefore, in this study, we assume that the mean and variance of $U$ are known. From the above assumption, $Y$ and $W$ are independent for the given $X$.

$$f_{Y,W|X}(y,w|x) = \frac{f_{Y,W,X}(y,w,x)}{f_X(x)} = \frac{f_{Y,W,U}(y,w,w-x)}{f_X(x)}$$

$$= \frac{f_{Y,X}(y,x)f_U(w-x)}{f_X(x)} = f_{Y|X}(y|x)f_{W|X}(w|x).$$

We use this conditional independence when we calculate the expectations.

*2.2. The Naive Estimator*

The naive estimator $\hat{\boldsymbol{\beta}}^{(N)} = (\hat{\beta}_0^{(N)}, \hat{\beta}_1^{(N)})'$ for $\boldsymbol{\beta} = (\beta_0, \beta_1)'$ is defined as the solution of the equation

$$S_n(\hat{\boldsymbol{\beta}}^{(N)}|\mathcal{X}) = \mathbf{0}_2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \tag{2}$$

where

$$S_n(\boldsymbol{b}|\mathcal{X}) = \frac{1}{n}\sum_{i=1}^{n}\{Y_i - \exp(b_0 + b_1 W_i)\}(1, W_i)'$$

is a function of indeterminant $\boldsymbol{b} = (b_0, b_1)'$ given $\mathcal{X} = (X_1, \ldots, X_n)'$. The naive estimator can be interpreted as the maximum likelihood estimator if we wrongly assume that $Y|W \sim Po(\exp(\beta_0 + \beta_1 W))$ because (2) is the log-likelihood equation for $Y|W \sim Po(\exp(\beta_0 + \beta_1 W))$. The correct distribution of $Y|W$ is

$$f_{Y|W}(y|w) = \frac{1}{f_W(w)} \int_{supp(f_U)} f_{Y|W,U}(y|w,u)f_U(u)f_X(w-u)du$$

$$= \frac{1}{f_W(w)} \int_{supp(f_U)} f_{Y|X}(y|w-u)f_U(u)f_X(w-u)du$$

$$= \frac{1}{f_W(w)} \int_{supp(f_U)} Po(\exp(\beta_0 + \beta_1(w-u)))f_U(u)f_X(w-u)du$$

assuming that $U$ is independent of $(X, Y|X)$. The right-hand side must be different from $Po(\exp(\beta_0 + \beta_1 W))$ in general. If one ignores the error $U$ and fits the likelihood estimation using $W$ instead of $X$, a biased estimator is obtained. In fact, by the law of large numbers, we have

$$S_n(\hat{\boldsymbol{\beta}}^{(N)}|\mathcal{X}) = \frac{1}{n}\sum_{i=1}^{n}\{Y_i - \exp(\hat{\beta}_0^{(N)} + \hat{\beta}_1^{(N)}W_i)\}(1, W_i)'$$

$$\xrightarrow{p} \mathbf{E}_{X,W}[\mathbf{E}_{Y|(X,W)}[\{Y - \exp(\hat{\beta}_0^{(N)} + \hat{\beta}_1^{(N)}W)\}(1, W)']].$$

Thus, the naive estimator converges to $\boldsymbol{b} = (b_0, b_1)'$ which is the solution of the estimating equation

$$\mathbf{E}_{X,W}[\mathbf{E}_{Y|(X,W)}[\{Y - \exp(\hat{\beta}_0^{(N)} + \hat{\beta}_1^{(N)}W)\}(1, W)']] = \mathbf{0}_2. \tag{3}$$

Equation (3) implies that for a given $\mathcal{X}$

$$\hat{\boldsymbol{\beta}}^{(N)} \xrightarrow{p} \boldsymbol{b} \neq \boldsymbol{\beta}.$$

The solution $\boldsymbol{b}$ of the estimating equation is generally different from $\boldsymbol{\beta}$.

## 3. Properties of the Naive Estimator

In this section, we consider the requirements for the existence of the naive estimator and derive its asymptotic bias and asymptotic MSE assuming that the explanatory variable and measurement error are generally distributed. In addition, we introduce application examples for a Gamma explanatory variable with a normal error or a Gamma error.

### 3.1. The Existence of the Naive Estimator

The naive estimator does not always exist for general random variables $X$ and $U$. Thus, we assume the existence of the expectation

$$\mathbf{E}_{X,Y,W}[\{Y - \exp(b_0 + b_1 W)\}(1, W)']$$

as a requirement for the existence of the naive estimator. Consequently, the following four expectations should exist.

$$\begin{cases} \mathbf{E}[Y] & = \mathbf{E}_X[\mathbf{E}[Y|X]] = \mathbf{E}_X[\exp(\beta_0 + \beta_1 X)] = e^{\beta_0} M_X(\beta_1), \\ \mathbf{E}[\exp(b_0 + b_1 W)] & = e^{b_0} \mathbf{E}[e^{b_1 X + b_1 U}] = e^{b_0} M_X(b_1) M_U(b_1), \\ \mathbf{E}[YW] & = \mathbf{E}_X[\mathbf{E}[Y|X]\mathbf{E}[W|X]] = \mathbf{E}_X[(X + \mathbf{E}[U]) \exp(\beta_0 + \beta_1 X)] \\ & = e^{\beta_0} \mathbf{E}[U] M_X(\beta_1) + e^{\beta_0} \mathbf{E}[X e^{\beta_1 X}] \\ & = e^{\beta_0} \mathbf{E}[U] M_X(\beta_1) + e^{\beta_0} \nabla M_X(\beta_1), \\ \mathbf{E}[W \exp(b_0 + b_1 W)] & = \mathbf{E}_X[\mathbf{E}_U[(X + U) \exp(b_0 + b_1 X + b_1 U)]] \\ & = e^{b_0} \mathbf{E}[X e^{b_1 X}] M_U(b_1) + e^{b_0} \mathbf{E}[U e^{b_1 U}] M_X(b_1) \\ & = e^{b_0} M_U(b_1) \nabla M_X(b_1) + e^{b_0} M_X(b_1) \nabla M_U(b_1). \end{cases} \quad (4)$$

Therefore, these expectations require that $M_X(\beta_1), M_X(b_1), M_U(b_1)$ exist. This condition is the requirement for the existence of the naive estimator. Here, we assume the existence of

$$M_X(\beta_1), M_X(b_1), M_U(b_1) \quad (5)$$

for the distributions of $X$ and $U$.

### 3.2. Asymptotic Bias of the Naive Estimator

The naive estimator satisfies

$$\hat{\boldsymbol{\beta}}^{(N)} \quad \xrightarrow{p} \quad \boldsymbol{b}$$

and has an asymptotic bias for the true $\boldsymbol{\beta}$. Here, we derive the asymptotic bias under general conditions. From (3), we obtain two equations:

$$\begin{cases} \mathbf{E}[Y] & = \mathbf{E}[\exp(b_0 + b_1 W)], \\ \mathbf{E}[YW] & = \mathbf{E}[W \exp(b_0 + b_1 W)]. \end{cases} \quad (6)$$

From (4) with the above equalities, we have

$$e^{\beta_0} M_X(\beta_1) = e^{b_0} M_X(b_1) M_U(b_1),$$

$$e^{\beta_0} \mathbf{E}[U] M_X(\beta_1) + e^{\beta_0} \nabla M_X(\beta_1) = e^{b_0} (\nabla M_X(b_1)) M_U(b_1) + e^{b_0} (\nabla M_U(b_1)) M_X(b_1)$$

$$= e^{b_0} \nabla (M_X(b_1) M_U(b_1)) = e^{b_0} \nabla M_W(b_1).$$

Therefore, we use a transformation to obtain the following system of equations:

$$\begin{cases} b_0 & = \beta_0 + \log\left(\frac{M_X(\beta_1)}{M_W(b_1)}\right), \\ K'_W(b_1) & = \frac{1}{M_W(b_1)} \nabla M_W(b_1) = \mathbf{E}[U] + \frac{\nabla M_X(\beta_1)}{M_X(\beta_1)}, \end{cases} \quad (7)$$

where $K_W$ is the cumulant generating function of $W$. Thus, $\boldsymbol{b} = (b_0, b_1)'$ is determined by the solution of this system of equations. Therefore, the equation

$$K'_W(b_1) = \mathbf{E}[U] + \frac{\nabla M_X(\beta_1)}{M_X(\beta_1)}$$

should have a solution with respect to $b_1$. Here, we set

$$G(\beta_1, b_1) := K'_W(b_1) - \mathbf{E}[U] - K'_X(\beta_1).$$

We assume $G(\beta_1, b_1)$ has zero in $\mathbb{R}^2$ and satisfies

$$\frac{\partial G(\beta_1, b_1)}{\partial b_1} = K''_W(b_1) \neq 0.$$

$G$ is continuously differentiable because we assume the existence of (5). Then, by the theorem of implicit functions, there exists a unique $C^1$-class function $g$ that satisfies $b_1 = g(\beta_1)$ in the neighborhood of the zero of $G$. Using this expression, we write the asymptotic bias of the naive estimator as

$$\lim_{n \to \infty} \mathbf{E}[\hat{\beta}_0^{(N)} - \beta_0] = b_0 - \beta_0 = \log\left(\frac{M_X(\beta_1)}{M_W \circ g(\beta_1)}\right),$$

$$\lim_{n \to \infty} \mathbf{E}[\hat{\beta}_1^{(N)} - \beta_1] = b_1 - \beta_1 = g(\beta_1) - \beta_1.$$

We also derive the asymptotic MSE of the naive estimator. The MSE can be represented as the sum of the squared bias and variance. The asymptotic variance of the naive estimator is $0$ because the naive estimator is a consistent estimator of $\boldsymbol{b}$. Thus, we obtain the asymptotic MSE of the naive estimator as

$$\lim_{n \to \infty} \mathbf{E}[(\hat{\beta}_0^{(N)} - \beta_0)^2] = (b_0 - \beta_0)^2 = \left(\log\left(\frac{M_X(\beta_1)}{M_W \circ g(\beta_1)}\right)\right)^2,$$

$$\lim_{n \to \infty} \mathbf{E}[(\hat{\beta}_1^{(N)} - \beta_1)^2] = (b_1 - \beta_1)^2 = (g(\beta_1) - \beta_1)^2.$$

Therefore, the asymptotic bias is given by the following theorem assuming general distributions.

**Theorem 1.** *Let $Y|X \sim Po(\exp(\beta_0 + \beta_1 X))$. Assume that $W = X + U$ and $U$ is independent of $(X, Y|X)$. Assume the existence of $M_X(\beta_1), M_X(b_1), M_U(b_1)$. Let*

$$G(\beta_1, b_1) := K'_W(b_1) - \mathrm{E}[U] - K'_X(\beta_1).$$

*Assume the function $G$ has a zero in $\mathbb{R}^2$, namely there exist solutions with $G(\beta_1, b_1) = 0$, and satisfies*

$$\frac{\partial G(\beta_1, b_1)}{\partial b_1} = K''_W(b_1) \neq 0.$$

*Then, the asymptotic biases of the naive estimators $\hat{\beta}_0^{(N)}$ and $\hat{\beta}_1^{(N)}$ are given by*

$$\log\left(\frac{M_X(\beta_1)}{M_W \circ g(\beta_1)}\right) \quad and \quad g(\beta_1) - \beta_1$$

*respectively, where $g$ is a $C^1$-class function satisfying $b_1 = g(\beta_1)$ in the neighborhood of the zero of $G$. Furthermore, the asymptotic MSEs of the naive estimators $\hat{\beta}_0^{(N)}$ and $\hat{\beta}_1^{(N)}$ are given by their squared asymptotic biases.*

*3.3. Examples*

In this section, we present two type of examples. First, we assume that a Gamma explanatory variable with a normal error. Let

$$X \sim \Gamma(k, \lambda), \quad U \sim N(0, \sigma^2),$$

where $k > 0, \lambda > 0, 0 < \sigma^2 < \infty$. We apply the naive estimation under this condition. From the assumptions of Theorem 1, we assume the existence of

$$M_X(\beta_1), M_X(b_1) \text{ and } M_U(b_1).$$

Therefore, we obtain the parameter conditions

$$\lambda - \beta_1 > 0, \quad \lambda - b_1 > 0.$$

Next, we derive $\boldsymbol{b} = (b_0, b_1)'$. Under this condition, we obtain

$$G(\beta_1, b_1) = K_W'(b_1) - \mathbf{E}[U] - K_X'(\beta_1) = \frac{k}{\lambda - b_1} + \sigma^2 b_1 - \frac{k}{\lambda - \beta_1}.$$

Thus, the set of zeros of $G$ is

$$\left\{ (\beta_1, b_1) \in \mathbb{R}^2; \beta_1 = \frac{k + \lambda \sigma^2 (\lambda - b_1)}{k + \sigma^2 (\lambda - b_1) b_1} b_1 \right\}.$$

In addition,

$$\frac{\partial G(\beta_1, b_1)}{\partial b_1} = \frac{k}{(\lambda - b_1)^2} + \sigma^2 > 0.$$

Therefore, $G$ has a zero in $\mathbb{R}^2$ and satisfies $\frac{\partial G(\beta_1, b_1)}{\partial b_1} \neq 0$. From $G(\beta_1, b_1) = 0$, we obtain two implicit functions

$$b_1^{(1)} = \frac{(\lambda - \beta_1)\lambda\sigma^2 + k + \sqrt{s}}{2(\lambda - \beta_1)\sigma^2},$$

$$b_1^{(2)} = \frac{(\lambda - \beta_1)\lambda\sigma^2 + k - \sqrt{s}}{2(\lambda - \beta_1)\sigma^2},$$

where $s = (\lambda - \beta_1)^2 \lambda^2 \sigma^4 + 2(\lambda - \beta_1)(\lambda - 2\beta_1)\sigma^2 k + k^2 > 0$. Then, we obtain two expressions of $b_0$ corresponding to $b_1$.

$$b_0^{(1)} := \beta_0 + \log\left( \frac{M_X(\beta_1)}{M_W(b_1^{(1)})} \right)$$

$$= \beta_0 + k \log \frac{(\lambda - \beta_1)\lambda\sigma^2 - k - \sqrt{s}}{2(\lambda - \beta_1)^2 \sigma^2}$$

$$- \frac{(\lambda - \beta_1)^2 \lambda^2 \sigma^4 + 2(\lambda - \beta_1)^2 \sigma^2 k + k^2 + ((\lambda - \beta_1)\lambda\sigma^2 + k)\sqrt{s}}{4(\lambda - \beta_1)^2 \sigma^2},$$

$$b_0^{(2)} := \beta_0 + \log\left( \frac{M_X(\beta_1)}{M_W(b_1^{(2)})} \right)$$

$$= \beta_0 + k \log \frac{(\lambda - \beta_1)\lambda\sigma^2 - k + \sqrt{s}}{2(\lambda - \beta_1)^2 \sigma^2}$$

$$- \frac{(\lambda - \beta_1)^2 \lambda^2 \sigma^4 + 2(\lambda - \beta_1)^2 \sigma^2 k + k^2 - ((\lambda - \beta_1)\lambda\sigma^2 + k)\sqrt{s}}{4(\lambda - \beta_1)^2 \sigma^2}.$$

In addition,

$$s = ((\lambda - \beta_1)\lambda\sigma^2 - k)^2 + 4(\lambda - \beta_1)^2\sigma^2 k;$$

therefore, $s$ satisfies $\sqrt{s} > | (\lambda - \beta_1)\lambda\sigma^2 - k |$. From the antilogarithm condition, $\boldsymbol{b} = (b_0^{(2)}, b_1^{(2)})'$ is a solution of the system of Equation (6) in the range of $\mathbb{R}^2$. Thus, the asymptotic biases are given by

$$b_0 - \beta_0 = k\log\frac{(\lambda - \beta_1)\lambda\sigma^2 - k + \sqrt{s}}{2(\lambda - \beta_1)^2\sigma^2}$$

$$- \frac{(\lambda - \beta_1)^2\lambda^2\sigma^4 + 2(\lambda - \beta_1)^2\sigma^2 k + k^2 - ((\lambda - \beta_1)\lambda\sigma^2 + k)\sqrt{s}}{4(\lambda - \beta_1)^2\sigma^2},$$

$$b_1 - \beta_1 = \frac{\lambda}{2} - \beta_1 + \frac{k - \sqrt{s}}{2(\lambda - \beta_1)\sigma^2}.$$

Next, we present another example, Gamma explanatory variable with a Gamma error. Let

$$X \sim \Gamma(k_1, \lambda), \quad U \sim \Gamma(k_2, \lambda),$$

where $k_1 > 0, k_2 > 0, \lambda > 0$. We apply the naive estimation under this condition. From the assumptions of Theorem 1, we assume the existence of

$$M_X(\beta_1), M_X(b_1) \text{ and } M_U(b_1).$$

Therefore, we obtain the parameter conditions

$$\lambda - \beta_1 > 0, \quad \lambda - b_1 > 0.$$

Next, we derive $\boldsymbol{b} = (b_0, b_1)'$. Under this condition, we obtain

$$G(\beta_1, b_1) = \frac{k_1 + k_2}{\lambda - b_1} - \frac{k_1}{\lambda - \beta_1} - \frac{k_2}{\lambda}.$$

Thus, the set of zeros of $G$ is

$$\left\{ (\beta_1, b_1) \in \mathbb{R}^2; b_1 = \frac{k_1\lambda\beta_1}{k_1\lambda + k_2(\lambda - \beta_1)} \right\}.$$

In addition,

$$\frac{\partial G(\beta_1, b_1)}{\partial b_1} = \frac{k_1 + k_2}{(\lambda - b_1)^2} > 0.$$

Therefore, $G$ has a zero in $\mathbb{R}^2$ and satisfies $\frac{\partial G(\beta_1, b_1)}{\partial b_1} \neq 0$. From $G(\beta_1, b_1) = 0$, we obtain the implicit function

$$b_1 = \frac{k_1\lambda\beta_1}{k_1\lambda + k_2(\lambda - \beta_1)}.$$

Thus, by Theorem 1, the asymptotic biases are given by

$$b_0 - \beta_0 = -k_1\log(1 - \beta_1/\lambda) + (k_1 + k_2)\log(1 - b_1/\lambda),$$

$$b_1 - \beta_1 = -\frac{k_2(\lambda - \beta_1)\beta_1}{k_1\lambda + k_2(\lambda - \beta_1)}.$$

## 4. Corrected Naive Estimator

In this section, we propose a corrected naive estimator as a consistent estimator of $\boldsymbol{\beta}$ under general distributions and give application examples for a Gamma explanatory variable with a normal error or a Gamma error. From (7), we have the following system of equations:

$$\beta_0 = b_0 + \log\left(\frac{M_W(b_1)}{M_X(\beta_1)}\right),$$

$$G(\beta_1, b_1) = K'_W(b_1) - \mathbf{E}[U] - K'_X(\beta_1) = 0.$$

By solving this system of equations for $\beta_0, \beta_1$ and replacing $\boldsymbol{b} = (b_0, b_1)'$ with the naive estimator $\hat{\boldsymbol{\beta}}^{(N)} = (\hat{\beta}_0^{(N)}, \hat{\beta}_1^{(N)})'$, we obtain the consistent estimator of the true $\boldsymbol{\beta}$. Here,

$$\hat{\boldsymbol{\beta}}^{(N)} = \begin{pmatrix} \hat{\beta}_0^{(N)} \\ \hat{\beta}_1^{(N)} \end{pmatrix} \qquad \xrightarrow{p} \qquad \boldsymbol{b} = \begin{pmatrix} b_0 \\ b_1 \end{pmatrix}.$$

Therefore,

$$\hat{\boldsymbol{\beta}}^{(CN)} \qquad \xrightarrow{p} \qquad \boldsymbol{\beta}.$$

Thus, $\hat{\boldsymbol{\beta}}^{(CN)}$ is a consistent estimator of $\boldsymbol{\beta}$. If $G$ has zero in $\mathbb{R}^2$ and satisfies

$$\frac{\partial G(\beta_1, b_1)}{\partial \beta_1} = -K''_X(\beta_1) \neq 0,$$

then, by the theorem of implicit functions, there exists a unique $C^1$-class function $h$ that satisfies $\beta_1 = h(b_1)$ in the neighborhood of the zero of $G$. We note that $h$ is the inverse function of $g$ in Theorem 1. We propose a corrected naive estimator that is the consistent estimator of the true $\boldsymbol{\beta}$ as follows.

**Theorem 2.** *Let $Y|X \sim Po(\exp(\beta_0 + \beta_1 X))$. Assume that $W = X + U$ and $U$ is independent of $(X, Y|X)$. Assume the existence of $M_X(\beta_1), M_X(b_1), M_U(b_1)$. Let*

$$G(\beta_1, b_1) := K'_W(b_1) - \boldsymbol{E}[U] - K'_X(\beta_1).$$

*Assume $G$ has zero in $\mathbb{R}^2$ and satisfies*

$$\frac{\partial G(\beta_1, b_1)}{\partial \beta_1} = -K''_X(\beta_1) \neq 0.$$

*Then, the corrected naive estimator $\hat{\boldsymbol{\beta}}^{(CN)} = (\hat{\beta}_0^{(CN)}, \hat{\beta}_1^{(CN)})'$, which corrects the bias of the naive estimator $\hat{\boldsymbol{\beta}}^{(N)} = (\hat{\beta}_0^{(N)}, \hat{\beta}_1^{(N)})'$, is given by*

$$\hat{\beta}_0^{(CN)} = \hat{\beta}_0^{(N)} + \log\left(\frac{M_W(\hat{\beta}_1^{(N)})}{M_X(\hat{\beta}_1^{(CN)})}\right),$$

$$\hat{\beta}_1^{(CN)} = h(\hat{\beta}_1^{(N)}),$$

*where $h$ is a $C^1$-class function satisfying $\beta_1 = h(b_1)$ in the neighborhood of the zero of $G$. Furthermore, the corrected naive estimator is a consistent estimator of $\boldsymbol{\beta}$.*

**Example 1.** *We derive the corrected naive estimator assuming*

$$X \sim \Gamma(k, \lambda), U \sim N(0, \sigma^2).$$

*We obtain*

$$G(\beta_1, b_1) = \frac{k}{\lambda - b_1} + \sigma^2 b_1 - \frac{k}{\lambda - \beta_1},$$

$$\frac{\partial G(\beta_1, b_1)}{\partial \beta_1} = -\frac{k}{(\lambda - \beta_1)^2} < 0.$$

*G has zero in $\mathbb{R}^2$ and satisfies $\frac{\partial G(\beta_1, b_1)}{\partial \beta_1} \neq 0$. From $G(\beta_1, b_1) = 0$, we obtain the implicit function*

$$\beta_1 = \frac{\sigma^2 \lambda b_1^2 - (k + \lambda^2 \sigma^2) b_1}{\sigma^2 b_1^2 - \lambda \sigma^2 b_1 - k} = h(b_1).$$

*Thus, by Theorem 2, the corrected naive estimator is given by*

$$\hat{\beta}_0^{(CN)} = \hat{\beta}_0^{(N)} + \log\left(\frac{M_W(\hat{\beta}_1^{(N)})}{M_X(\hat{\beta}_1^{(CN)})}\right)$$

$$= \hat{\beta}_0^{(N)} + \frac{1}{2}\hat{\beta}_1^{(N)2}\sigma^2 + k\log(1 - \hat{\beta}_1^{(CN)}/\lambda) - k\log(1 - \hat{\beta}_1^{(N)}/\lambda),$$

$$\hat{\beta}_1^{(CN)} = h(\hat{\beta}_1^{(N)}) = \frac{\lambda\sigma^2\hat{\beta}_1^{(N)2} - (k + \lambda^2\sigma^2)\hat{\beta}_1^{(N)}}{\sigma^2\hat{\beta}_1^{(N)2} - \lambda\sigma^2\hat{\beta}_1^{(N)} - k}.$$

**Example 2.** *We derive the corrected naive estimator assuming*

$$X \sim \Gamma(k_1, \lambda), U \sim \Gamma(k_2, \lambda).$$

*We obtain*

$$G(\beta_1, b_1) = \frac{k_1 + k_2}{\lambda - b_1} - \frac{k_1}{\lambda - \beta_1} - \frac{k_2}{\lambda},$$

$$\frac{\partial G(\beta_1, b_1)}{\partial \beta_1} = -\frac{k_1}{(\lambda - \beta_1)^2} < 0.$$

*G has zero in $\mathbb{R}^2$ and satisfies $\frac{\partial G(\beta_1, b_1)}{\partial \beta_1} \neq 0$. From $G(\beta_1, b_1) = 0$, we obtain the implicit function*

$$\beta_1 = \frac{(k_1 + k_2)b_1\lambda}{k_1\lambda + k_2 b_1} = h(b_1).$$

*Thus, by Theorem 2, the corrected naive estimator is given by*

$$\hat{\beta}_0^{(CN)} = \hat{\beta}_0^{(N)} + \log\left(\frac{M_W(\hat{\beta}_1^{(N)})}{M_X(\hat{\beta}_1^{(CN)})}\right)$$

$$= \hat{\beta}_0^{(N)} + k_1\log(1 - \hat{\beta}_1^{(CN)}/\lambda) - (k_1 + k_2)\log(1 - \hat{\beta}_1^{(N)}/\lambda),$$

$$\hat{\beta}_1^{(CN)} = h(\hat{\beta}_1^{(N)}) = \frac{(k_1 + k_2)\hat{\beta}_1^{(N)}\lambda}{k_1\lambda + k_2\hat{\beta}_1^{(N)}}.$$

## 5. Simulation Studies

In this section, we present simulation studies that compare the performance of the naive estimator and corrected naive estimator. We denote the sample size by $n$ and the number of simulations by MC. We calculate the estimated bias for $\hat{\boldsymbol{\beta}}^{(N)}$ and $\hat{\boldsymbol{\beta}}^{(CN)}$ as follows:

$$\widehat{\text{BIAS}(\hat{\boldsymbol{\beta}}^{(N)})} = \frac{1}{MC}\sum_{i=1}^{MC}\hat{\boldsymbol{\beta}}_i^{(N)} - \boldsymbol{\beta},$$

$$\widehat{\text{BIAS}(\hat{\boldsymbol{\beta}}^{(CN)})} = \frac{1}{MC}\sum_{i=1}^{MC}\hat{\boldsymbol{\beta}}_i^{(CN)} - \boldsymbol{\beta},$$

where $\hat{\beta}_i^{(N)}$ and $\hat{\beta}_i^{(CN)}$ represent the naive estimator and corrected naive estimator in the $i$th time simulation, respectively. Similarly, we calculate the estimated MSE matrix for $\hat{\beta}^{(N)}$ and $\hat{\beta}^{(CN)}$ as follows:

$$\widehat{\mathrm{MSE}(\hat{\boldsymbol{\beta}}^{(N)})} = \frac{1}{MC} \sum_{i=1}^{MC} (\hat{\boldsymbol{\beta}}_i^{(N)} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}}_i^{(N)} - \boldsymbol{\beta})',$$

$$\widehat{\mathrm{MSE}(\hat{\boldsymbol{\beta}}^{(CN)})} = \frac{1}{MC} \sum_{i=1}^{MC} (\hat{\boldsymbol{\beta}}_i^{(CN)} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}}_i^{(CN)} - \boldsymbol{\beta})'.$$

*5.1. Case 1*

We assume $X \sim \Gamma(k, \lambda), U \sim N(0, \sigma^2)$. Let $\beta_0 = 0.2, \beta_1 = 0.3, k = 2, \lambda = 1.2$, $n = 500, MC = 1000$. We perform simulations with $\sigma^2 = 0.05, 0.5, 2$. Note that we assume that the true value of $\sigma^2$ is known. We estimate $k, \lambda$ in the formula of the corrected naive estimator by the moment method in terms of $W$ because the value of $X$ cannot be directly observed.

$$\hat{k} = \left( \frac{1}{n} \sum_{i=1}^{n} w_i \right) \hat{\lambda},$$

$$\hat{\lambda} = \frac{\frac{1}{n} \sum_{i=1}^{n} w_i}{\frac{1}{n} \sum_{i=1}^{n} (w_i - \bar{w})^2 - \sigma^2},$$

where $w_i$ $(i = 1, \ldots, n)$ is the samples of $W$.

Table 1 shows the estimated bias of the true $\boldsymbol{\beta}$. Asy.Bias $\hat{\beta}_0$ and Asy.Bias $\hat{\beta}_1$ denote the theoretical asymptotic biases of $\hat{\beta}_0^{(N)}$ and $\hat{\beta}_1^{(N)}$, respectively, given in Theorem 1. The bias correction of the naive estimator is performed by the corrected naive estimator. With increasing $\sigma^2$, the bias of the naive estimator increases. However, the bias of the corrected naive estimator is small for large $\sigma^2$.

**Table 1.** Estimated bias of a Gamma distribution with a Normal error.

|  |  | Asy.Bias $\hat{\beta}_0$ | $\widehat{\mathrm{BIAS}(\hat{\beta}_0)}$ | Asy.Bias $\hat{\beta}_1$ | $\widehat{\mathrm{BIAS}(\hat{\beta}_1)}$ |
|---|---|---|---|---|---|
| $\sigma^2 = 0.05$ | Naive | 0.01111 | 0.01139 | −0.005993 | −0.007199 |
|  | CN | 0 | 0.00003532 | 0 | 0.0002603 |
| $\sigma^2 = 0.5$ | Naive | 0.09912 | 0.1025 | −0.05297 | −0.05582 |
|  | CN | 0 | 0.007817 | 0 | 0.0007142 |
| $\sigma^2 = 2$ | Naive | 0.2757 | 0.2774 | −0.1454 | −0.1472 |
|  | CN | 0 | −0.009493 | 0 | 0.002736 |

Table 2 shows the estimated MSE of the true $\boldsymbol{\beta}$. Asy.MSE $\hat{\beta}_0$ and Asy.MSE $\hat{\beta}_1$ denote the theoretical asymptotic MSEs of $\hat{\beta}_0^{(N)}$ and $\hat{\beta}_1^{(N)}$, respectively, given in Theorem 1. The MSE of the corrected naive estimator is smaller than that of the naive estimator in all cases.

**Table 2.** Estimated MSE of a Gamma distribution with a normal error.

|  |  | Asy.MSE $\hat{\beta}_0$ | $\widehat{\mathrm{MSE}(\hat{\beta}_0)}$ | Asy.MSE $\hat{\beta}_1$ | $\widehat{\mathrm{MSE}(\hat{\beta}_1)}$ |
|---|---|---|---|---|---|
| $\sigma^2 = 0.05$ | Naive | 0.0001235 | 0.003003 | 0.00003592 | 0.0004536 |
|  | CN | 0 | 0.002920 | 0 | 0.0004254 |
| $\sigma^2 = 0.5$ | Naive | 0.009824 | 0.01362 | 0.002806 | 0.003508 |
|  | CN | 0 | 0.003806 | 0 | 0.0006354 |
| $\sigma^2 = 2$ | Naive | 0.07600 | 0.08124 | 0.02115 | 0.02214 |
|  | CN | 0 | 0.01021 | 0 | 0.002160 |

*5.2. Case 2*

We assume $X \sim \Gamma(k_1, \lambda), U \sim \Gamma(k_2, \lambda)$. Let $\beta_0 = 0.2, \beta_1 = 0.3, k_1 = 2, \lambda = 1.2, n = 500,$ $MC = 1000$. We perform simulations with $k_2 = 0.072, 0.72, 2.88$. Similarly, we assume that the true value of $k_2$ is known. We estimate $k_1, \lambda$ in the formula of the corrected naive estimator by the moment method in terms of $W$ because the value of $X$ cannot be directly observed.

$$\hat{k}_1 = \left( \frac{1}{n} \sum_{i=1}^{n} w_i \right) \hat{\lambda} - k_2,$$

$$\hat{\lambda} = \frac{\frac{1}{n} \sum_{i=1}^{n} w_i}{\frac{1}{n} \sum_{i=1}^{n} (w_i - \bar{w})^2},$$

where $w_i$ $(i = 1, \ldots, n)$ is the samples of $W$.

Table 3 shows the estimated bias of the true $\beta$. Similarly, the bias correction of the naive estimator is performed by the corrected naive estimator. The bias of the corrected naive estimator is small when the variance of the error is large. Table 4 shows the estimated MSE of the true $\beta$. The MSE of the corrected naive estimator is also smaller than that of the naive estimator.

**Table 3.** Estimated bias of a Gamma distribution with a Gamma error.

|  |  | **Asy.Bias $\hat{\beta}_0$** | $\widehat{\text{BIAS}(\hat{\beta}_0)}$ | **Asy.Bias $\hat{\beta}_1$** | $\widehat{\text{BIAS}(\hat{\beta}_1)}$ |
|---|---|---|---|---|---|
| $k_2 = 0.072$ | Naive | −0.002634 | −0.005415 | −0.007887 | −0.008874 |
|  | CN | 0 | −0.0006636 | 0 | 0.0002777 |
| $k_2 = 0.72$ | Naive | −0.02090 | −0.01725 | −0.06378 | −0.06475 |
|  | CN | 0 | -0.0002963 | 0 | −0.003184 |
| $k_2 = 2.88$ | Naive | −0.04953 | −0.05439 | −0.1558 | −0.1569 |
|  | CN | 0 | 0.002954 | 0 | −0.003224 |

**Table 4.** Estimated MSE of a Gamma distribution with a Gamma error.

|  |  | **Asy.MSE $\hat{\beta}_0$** | $\widehat{\text{MSE}(\hat{\beta}_0)}$ | **Asy.MSE $\hat{\beta}_1$** | $\widehat{\text{MSE}(\hat{\beta}_1)}$ |
|---|---|---|---|---|---|
| $k_2 = 0.072$ | Naive | 0.08533 | 0.003109 | 0.000006940 | 0.0005384 |
|  | CN | 0 | 0.003074 | 0 | 0.0004743 |
| $k_2 = 0.72$ | Naive | 0.05580 | 0.005320 | 0.0004368 | 0.004894 |
|  | CN | 0 | 0.004457 | 0 | 0.0008818 |
| $k_2 = 2.88$ | Naive | 0.02080 | 0.01147 | 0.002453 | 0.02553 |
|  | CN | 0 | 0.007401 | 0 | 0.001963 |

## 6. Real Data Analysis

In this section, we apply the naive and corrected naive estimators to real data in two cases. First, we consider football data provided by Understat (2014). In this work, we focus on Goals and expected Goals (xG) in data on $N$ = 24,580 matches over 6 seasons between 2014–2015 and 2019–2020 from the Serie A, the Bundesliga, La Liga, the English Premier League, Ligue 1, and the Russian Premier League. Detail, such as the types and descriptions of the features, used in this section are provided in Table 5.

**Table 5.** Details of the variables.

| Features | Type | Description |
|---|---|---|
| Goals | counting | number of goals scored in the match |
| xG | continuous | performance metric used to evaluate football team and player performance |

We use goals as an objective variable $Y$ and xG as an explanatory variable $X$ and assume $Y|X \sim Po(\exp(\beta_0 + \beta_1 X))$ as the true model. Thus, this Poisson regression model refers to the extent to which expected goals (xG) explains (true) goals. We assume that the true parameter $\beta$ is obtained by the estimate from all $N$ data.

As a diagnostic technique, we calculate a measure of goodness-of-fit to verify that the dataset follows a Poisson regression model. Table 6 shows estimates of $\phi$ and $R_{McF}$ (McFadden 1974), where $R_{McF}$ is the ratio of the log-likelihood estimate to the initial log-likelihood. $\phi = \mathbf{V}[Y|X]/\mathbf{E}[Y|X]$ is an overdispersion parameter. We may consider that overdispersion is not observed because $\phi = 1$ equates to the standard Poisson regression model. The estimated value of $\beta$ is $(-0.5225, 0.5308)'$. Thus, we use this estimate as a true value. We assume $X$ (xG) $\sim \Gamma(k_1, \lambda)$ and obtain estimates of $k_1, \lambda$ as $k_1 = 2.425$, $\lambda = 1.851$ (see Figure 1).
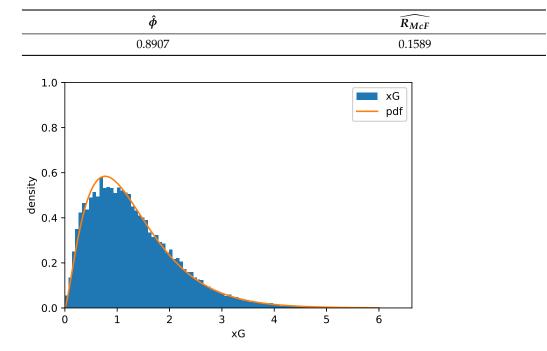
**Table 6.** Estimates of $\phi$ and $R_{McF}$.

| $\hat{\phi}$ | $\widehat{R_{McF}}$ |
|---|---|
| 0.8907 | 0.1589 |



**Figure 1.** Distribution of xG.

Expected goals (xG) is a performance metric used to represent the probability of a scoring opportunity that may result in a goal. xG is typically calculated from shot data. The measurer assigns a probability of scoring to a given shot and calculates the sum of the probabilities over a single game as xG. Observation error may occur in subjective evaluations. We can consider the situation that a high scorer happened to rate. Thus, we assume that $X$ includes a stochastic error $U$ given as

$$W = X + U.$$

Because $W$ must be a positive value, we choose a positive error by $U \sim \Gamma(k_2, \lambda)$ with $k_2 = k_1/10, k_1/3, k_1$. We sample 1000 random samples from among all $N$ samples to obtain the values of the estimates of $\beta$s. We repeat the estimations $MC = 10{,}000$ times to obtain the Monte Carlo mean of $\beta$s. The bias is calculated by the difference between the Monte Carlo mean and the true value.

Table 7 shows the estimated bias calculated by 10,000 simulations. The estimated bias of the corrected naive estimator is smaller than that of the naive estimator in all cases.

**Table 7.** Estimated bias and asymptotic bias in football data.

|  |  | Asy.Bias $\hat{\beta}_0$ | $\widehat{\text{BIAS}}(\hat{\beta}_0)$ | Asy.Bias $\hat{\beta}_1$ | $\widehat{\text{BIAS}}(\hat{\beta}_1)$ |
|---|---|---|---|---|---|
| $k_2 = k_1/10$ | Naive | −0.01148 | −0.01337 | −0.03534 | −0.03471 |
| | CN | 0 | −0.001804 | 0 | 0.0006200 |
| $k_2 = k_1/3$ | Naive | −0.03263 | −0.02383 | −0.1020 | −0.1067 |
| | CN | 0 | 0.008176 | 0 | −0.005575 |
| $k_2 = k_1$ | Naive | −0.06889 | −0.04692 | −0.2210 | −0.2291 |
| | CN | 0 | 0.01871 | 0 | -0.01215 |

Next, we apply the naive and corrected naive estimators to financial data based on data collected in the FinAccess survey conducted in 2019, provided by Kenya National Bureau of Statistics (2019). In this study, we focus on the values labelled as finhealthscore and Normalized Household weights, with a sample size of $N = 8669$. Details of the features used in this section, such as their types and descriptions, are provided in Table 8.

**Table 8.** Details of the variables.

| Features | Type | Description |
|---|---|---|
| finhealthscore | counting | Score of financial health for households |
| Normalized Household weights | continuous | Weighted and normalized households |

We use finhealthscore as an objective variable $Y$ and normalized household weights as an explanatory variable $X$ and assume $Y|X \sim Po(\exp(\beta_0 + \beta_1 X))$ as the true model. We further assume that the true parameter $\boldsymbol{\beta}$ is obtained by the estimate from all $N$ data.

As a diagnostic technique, we calculate a measure of goodness-of-fit to verify that the dataset follows a Poisson regression model. Table 9 shows estimates of $\phi$ and $R_{McF}$ (McFadden 1974). Overdispersion tends to occur to some extent in this Poisson regression model because the estimate of $\phi$ is greater than 1. The estimated value of $\boldsymbol{\beta}$ is $(1.0442, 0.1568)'$. As in the previous example, we regard the estimate as a true value. We assume $X \sim \Gamma(k_1, \lambda)$ and obtain estimates of $k_1, \lambda$ as $k_1 = 2.0746$, $\lambda = 2.0746$ (see Figure 2).

**Table 9.** Estimates of $\phi$ and $R_{McF}$.

| $\hat{\phi}$ | $\widehat{R_{McF}}$ |
|---|---|
| 1.4360 | 0.4478 |

According to Kenya National Bureau of Statistics (2019), the data from the FinAccess survey were weighted and adjusted for non-responses to obtain a representative dataset at the national and county level. Thus, we may consider the situation that $X$ exhibits a stochastic error $U$ as

$$W = X + U.$$

We assume a positive error by $U \sim \Gamma(k_2, \lambda)$ with $k_2 = k_1/10, k_1/3, k_1$ because the distribution of normalized household weights is positive. We sample random 1000 samples from among all $N$ samples to obtain the values of the estimates of $\beta$s. We repeat the estimations over $MC = 10,000$ iterations to obtain the Monte Carlo mean of $\beta$s. The bias is calculated by the difference between the Monte Carlo mean and the true value.
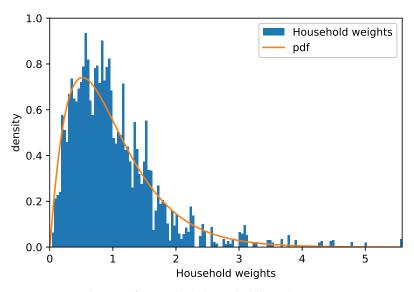
**Figure 2.** Distribution of normalized household weights.

Table 10 shows estimated bias calculated by 10,000 simulations. The estimated bias of the corrected naive estimator is smaller than that of the naive estimator in all cases.

**Table 10.** Estimated bias and asymptotic bias in financial data.

|  |  | **Asy.Bias $\hat{\beta}_0$** | $\widehat{\textbf{BIAS}(\hat{\beta}_0)}$ | **Asy.Bias $\hat{\beta}_1$** | $\widehat{\textbf{BIAS}(\hat{\beta}_1)}$ |
|---|---|---|---|---|---|
| $k_2 = k_1/10$ | Naive | −0.0005704 | −0.002225 | −0.01327 | −0.01207 |
|  | CN | 0 | −0.001628 | 0 | 0.001275 |
| $k_2 = k_1/3$ | Naive | −0.001581 | −0.004088 | −0.03694 | −0.03522 |
|  | CN | 0 | -0.002404 | 0 | 0.002119 |
| $k_2 = k_1$ | Naive | −0.003204 | −0.008314 | −0.07534 | −0.07283 |
|  | CN | 0 | −0.004744 | 0 | 0.004338 |

## 7. Discussion

In this study, we have proposed a corrected naive estimator as a consistent estimator for a Poisson regression model with a measurement error. Although Kukush et al. (2004) showed that the naive estimator has an asymptotic bias, the authors did not provide a method to correct this bias. Therefore, we developed an approach to estimate a Poisson regression model with an error. In contrast, the authors of Kukush et al. (2004) also proposed a corrected score estimator and a structural quasi-score estimator for a Poisson regression model with an error. These estimators are score-based and consistent for unknown parameters. Hence, a generalization of these estimators should be considered in future research. In addition, the model considered in the present work is restricted in the univariate case. Extending the explanatory variable to the multivariate case also remains a challenge of note.

**Author Contributions:** K.W. mainly worked this study supported by the second named author. K.W.: Derivation of the formulae, Proof of propositions, Application to the specific problems, Conduct the simulation study, Real data analysis, Coding of the programs. T.K.: Basic idea, Theoretical advice of the proof, Advice at each step, Whole checking. All authors have read and agreed to the published version of the manuscript.

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** Data were obtained from https://understat.com/ and https://knbs.or.ke (accessed on 11 February 2023).

## References

Box, George Edward Pelham. 1963. The Effects of Errors in the Factor Levels and Experimental Design. *Technometrics* 5: 247–62. [CrossRef]

Burr, Deborah. 1988. On Errors-in-Variables in Binary Regression-Berkson Case. *Journal of the American Statistical Association* 83: 739–43.

Geary, ROBERT C. 1953. Non-Linear Functional Relationship between Two Variables When One Variable is Controlled. *Journal of the American Statistical Association* 48: 94–103. [CrossRef]

Guo, Jie Q., and Tong Li. 2002. Poisson regression models with errors-in-variables:implication and treatment. *Journal of Statistical Planning and Inference* 104: 391–401. [CrossRef]

Huwang, Longcheen, and YH Steve Huang. 2000. On error-in-variables in polynomial regression-Berkson case. *Statistica Sinica* 10: 923–36.

Jiang, Fei, and Yanyuan Ma. 2020. Poisson Regression with Error Corrupted High Dimensional Features. *Statistica Sinica* 32: 2023–46. [CrossRef]

Kenya National Bureau of Statistics (KNBS). 2019. Available online: https://knbs.or.ke (accessed on 11 February 2023).

Kukush, Alexander, and Hans Schneeweiss. 2000. A Comparison of Asymptotic Covariance Matrices of Adjusted Least Squares and Structural Least Squares in Error Ridden Polynomial Regression. *Sonderforschungsbereich* 386: Paper 218. [CrossRef]

Kukush, Alexander, Hans Schneeweis, and Roland Wolf. 2004. Three Estimators for the Poisson Regression Model with Measurement Errors. *Statistical Papers* 45: 351–68. [CrossRef]

McFadden, Daniel. 1974. Conditional logit analysis of qualitative choice behavior. *Computations in Statistics-Theory and Methods* 47: 105–42.

Nakamura, Tsuyoshi. 1990. Corrected score function for errors-in-variables models: Methodology and application to generalized linear models. *Biometrika* 77: 127–37. [CrossRef]

Shklyar, Schneeweiss, and Hans Schneeweiss. 2005. A comparison of asymptotic covariance matrices of three consistent estimators in the Poisson regression model with measurement errors. *Journal of Multivariate Analysis* 94: 250–70. [CrossRef]

Understat. 2014. Available online: https://understat.com/ (accessed on 11 February 2023).