

Application of Tobit Regression on Household Expenditure on Egg and Milk Consumption in Bengkulu City

Claudia Nevani^{1*}, Sigit Nugroho², Winalia Agwil³

^{1,2,3} S1 Statistics Study Program, Bengkulu University, Bengkulu

* Corresponding Author: claudianevani@gmail.com

Article Info

Article History:

Received: 18 02 2025

Revised: 23 04 2025

Accepted: 23 04 2025

Available Online: 24 04 2025

Key Words:

Censored Data

Egg and Milk Consumption

Household Expenditure

Tobit Regression

Abstract

Regression analysis is a statistical method used to examine the functional relationship between two or more independent variables and a dependent variable. One of the regression methods designed to handle censored data or data with significant zero values is Tobit regression. This study aims to model household expenditures on egg and milk consumption in Bengkulu City using Tobit regression and to identify the factors influencing these expenditures. The data were obtained from the 2022 National Socioeconomic Survey, with a total sample of 1,170 households. The Tobit regression model was chosen because most household expenditure data had zero values, indicating censored data characteristics. This study identified several factors affecting expenditures on egg and milk consumption, such as the household head's education level, the number of household members, and the household head's employment sector. The results showed that the education level of the household head (elementary, junior high, and high school), the number of household members, and the household head's employment in agriculture and trade sectors had significant impacts on household expenditures for egg and milk consumption. The education level of the household head and their employment sector had a negative relationship, while the number of household members showed a positive relationship with these expenditures. The Tobit regression model successfully modeled household expenditures with adequate accuracy, as indicated by a Mean Absolute Percentage Error (MAAPE) of 1.38%.

1. INTRODUCTION

Regression analysis is a statistical method used to understand the relationship between independent variables and dependent variables. If the dependent variable is continuous, linear regression is used, while for categorical dependent variables, logistic regression is appropriate. When the dependent variable includes censored data (certain values are not fully observed), Tobit regression is the appropriate method. Tobit regression uses the Maximum Likelihood Estimation (MLE) approach to produce consistent and efficient estimates [1].

First introduced by [3], Tobit regression is used to analyze data such as household expenditures. For example, expenditures for egg and milk consumption are indicators of household protein adequacy. This data is often censored, with many households not allocating a budget for eggs and milk, as recorded in the National Socio-Economic Survey (SUSENAS) by BPS.

Previous studies have shown the success of Tobit regression in various contexts, such as student clothing expenditures [4] and household clean water consumption. Continuing this relevance, the study "Application of Tobit Regression on Household Egg and Milk Consumption Expenditures in Bengkulu City" aims to identify factors that influence budget allocation for the consumption of these protein sources.

2 METHOD

2.1 Censored Data

Censored data refers to data with incomplete observations because the value of a particular variable is limited to a certain range, and cannot be fully observed. There are three types of censoring in Tobit regression [2]:

1. Right Censoring: The value of the dependent variable is unknown if it exceeds a certain limit.
2. Left Censoring: The value of the dependent variable is unknown if it falls below a certain limit.

3. Interval Censoring: The value of the dependent variable is known to fall within a certain interval, but the exact value is unknown.

2.2 Tobit Model

The Tobit model was introduced by [3] to handle censored data. The general equation of the Tobit model is as follows :

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, 2, \dots, n$$

with y_i^* as the latent (unobserved) variables, where x_i is the vector of independent variables for the i -th observation, $\boldsymbol{\beta}$ is a vector of regression coefficient parameters, ε_i is the error term, n is the total number of observations, and:

$$y_i = \begin{cases} y_i^*, & \text{jika } y^* > c \\ c, & \text{jika } y^* \leq c \end{cases}$$

This model uses the Maximum Likelihood Estimation (MLE) method to estimate parameters because censored data cannot be processed using classical linear regression.

2.3 Censored Normal Distribution

The censored normal distribution is used for the dependent variable in the Tobit model. Its probability density function is expressed as:

$$f(y_i | \mathbf{x}_i) = \begin{cases} \frac{1}{\sigma} \phi\left(\frac{y_i - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right), & \text{jika } y_i > 0 \\ 1 - \phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right), & \text{jika } y_i = 0 \end{cases}$$

2.4 Maximum Likelihood Estimator (MLE)

The appropriate parameter estimation method for Tobit regression is the Maximum Likelihood Estimator (MLE) because it is consistent and efficient, especially for large samples [5].

The likelihood function of the standard Tobit model is:

$$L = \prod_{y_i=0}^n \left[1 - \phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right] \prod_{y_i>0}^n \left[\frac{1}{\sigma} \phi\left(\frac{y_i - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right]$$

Log likelihood function:

$$\begin{aligned} \ln L &= \sum_{y_i=0}^n \ln \left[1 - \phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right] + \sum_{y_i>0}^n \ln \left[\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y_i - \mathbf{x}_i' \boldsymbol{\beta})^2} \right] \\ &= \sum_{y_i=0}^n \ln \left[1 - \phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \right] - \frac{n}{2} \ln \sigma^2 - \frac{n}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \sum_{y_i>0}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 \end{aligned}$$

First derivative with respect to $\boldsymbol{\beta}$ and:

$$\frac{\partial \ln L}{\partial \boldsymbol{\beta}} = -\frac{1}{\sigma} \sum_{y_i=0}^n \frac{\phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \mathbf{x}_i'}{1 - \phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right)} + \frac{1}{\sigma^2} \sum_{y_i>0}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 \mathbf{x}_i'$$

and

$$\frac{\partial \ln L}{\partial \sigma^2} = -\frac{1}{2\sigma^3} \sum_{y_i=0}^n \frac{(\mathbf{x}_i' \boldsymbol{\beta}) \phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) \mathbf{x}_i'}{1 - \phi\left(\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right)} - \frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{y_i>0}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2$$

Thus, the parameter estimation equations. $\hat{\beta}$ and $\hat{\sigma}^2$ can be written as follows:

$$\hat{\beta} = (X'X)^{-1}X'Y - \sigma(X'X)^{-1}\bar{X}\bar{Y}$$

and

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{y_i > 0}^n (y_i - x_i'\beta)y_i$$

2.5 Assumption Testing

Several assumption tests are applied in the Tobit regression analysis:

1. Normality Test: Performed to check whether the residuals are normally distributed, for example using the Jarque-Bera test [6].

Hypotheses:

H_0 : Residuals are normally distributed.

H_1 : Residuals are not normally distributed.

Test statistics:

$$JB = n \left(\frac{S^2}{6} + \frac{(K-3)^2}{24} \right)$$

If the $p - value < 0.05$, then H_0 is rejected, indicating the residuals are not normally distributed. The JB statistic is compared to a chi-square distribution.

2. Multicollinearity Detection: Multicollinearity is detected by calculating the Variance Inflation Factor (VIF) VIF formula:

$$VIF_j = \frac{1}{1 - R_j^2}$$

A VIF value > 10 indicates high multicollinearity, R_j^2 is the coefficient of determination (R-squared) obtained from the regressing the j -th independent variable all other independent variables.

3. Heteroscedasticity Detection: The Breusch-Pagan test is used to detect inconsistencies in residual variances [7].

Hypotheses:

H_0 : No heteroscedasticity (homoscedastic residuals).

H_1 : Presence heteroscedasticity

Test statistics:

$$\phi = \frac{1}{2} (ESS)$$

If the $p - value < 0.05$, then H_0 is rejected, indicating heteroscedasticity, ESS stands for Explained Sum of Squares. The $p-value$ is compared to the Chi-square distribution.

2.6 Parameter Testing

Parameter testing in the Tobit model can be performed using the G test for simultaneous testing and the Wald Test for partial testing.

1. Simultaneous Test: This test uses the likelihood ratio method to evaluate whether the independent variables, collectively, have a significant influence on the dependent variable. The hypothesis being tested is:

Hypotheses:

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$ (all independent variables have no effect on the model)

H_1 : there is at least one $\beta_j \neq 0, j = 1, 2, 3, \dots, k$ (at least one independent variable has an effect on the model)

Test statistics:

$$G = -2 \ln \left(\frac{L(\hat{\omega})}{L(\hat{\Omega})} \right)$$

Reject H_0 if the test statistic exceeds the critical value ($G \geq \chi^2_{(\alpha; db=k)}$) or if the p-value is less than 0.05, indicating that at least one independent variables significantly influence in the model.

2. Partial Test: The Wald test tests assesses the effect of each independent variable individually.

Hypotheses:

$H_0: \beta_j = 0$ (variable X_j has no effect)

$H_1: \beta_j \neq 0, j = 1, 2, \dots, k$ (variable X_j has an effect)

Test statistics;

$$W = \left(\frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} \right)^2$$

Reject H_0 if the test statistic exceeds the critical value or if the p-value is less than 0.05, indicating that variable j significantly influences the model.

2.7 Mean Arctangent Absolute Percent Error (MAAPE)

MAAPE is an evaluation metric derived from MAPE designed to address the issue of infinite or undefined values when actual values are zero or near zero [8]. MAAPE is calculated using the formula:

$$MAAPE = \left(\frac{1}{n} \right) \sum_{i=1}^n \arctan \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100$$

A lower the MAAPE value, indicates better forecasting accuracy. According to Maricar (2019), the interpretation of MAAPE values is

Tabel 1. Range Nilai MAAPE

MAAPE (%)	Interpretation
< 10	Very accurate prediction
10 – 20	Good prediction
20 – 50	Reasonable prediction
> 50	Inaccurate prediction

MAAPE is a more stable alternative to MAPE, especially for data with low actual values.

This study uses secondary data from the results of SUSENAS in Bengkulu City in September 2022, with a sample of 1,170 households. The data collected includes household expenditure for egg and milk consumption, education level of the head of the household, number of household members, and employment of the head of the household.

The analysis is carried out through the following stages:

1. Tobit regression modeling.

The general Tobit regression model is an approach used in statistical analysis when the dependent variable is censored, meaning some values of the dependent variable are not fully observed due to being truncated at a certain threshold. In this model, it is assumed that there exists an underlying latent variable y_i^* that follows a standard linear regression model, represented as $y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i$.

2. Tobit regression assumption test:

- a. Normality (Jarque-Berra test).
- b. Multicollinearity (VIF).
- c. Homoscedasticity (Breusch-Pagan LM test).

3. Model parameter test:

- a. Simultaneous (Likelihood Ratio/G Test).
- b. Partial (Wald Test).

4. Model evaluation using MAAPE test.
5. Interpretation of tobit regression results

3 RESULT AND DISCUSSION

3.1 Tobit Regression Parameter Estimation

The parameter β_j in the model is estimated using the maximum likelihood method. The estimated value were obtained using R software and are presented in Appendix 3. The tobit regression estimation results for the variable influencing household expenditure on egg and milk consumption in Bengkulu City are shown in table 2:

Tabel 2. Parameter Estimation

Variable	$\hat{\beta}_j$	Standar Error	z-Value	P-Value	
Intercept	51.246,32	8.489,81	6.036	1.58×10^{-09}	***
DPD1	-29.474,67	6.800,47	-4.334	1.46×10^{-05}	***
DPD2	-21.660,52	7.269,15	-2.980	0.00288	**
DPD3	-14.753,86	6.432,38	-2.294	0.2181	*
yJART	9.433,46	1.297,61	7.270	3.60×10^{-15}	***
JA5T	571,19	2.980,93	0.192	0.84804	
DPK1	-14.967,72	7.461,24	-2.006	0.04485	*
DPK2	-21.677,86	14.437,11	-1.502	0.13322	
DPK3	-18.340,26	8.296,72	-2.211	0.02707	*
DPK4	-12.413,49	7.528,03	-1.649	0.09915	.
DPK5	-9.292,64	1.257,58	-0.739	0.45995	

Significance code: 0,001 = ***, 0,01 = **, 0,05 = *, 0,1 = .

the estimated Tobit regression model is as follows:

$$y^* = 51.246,32 - 29.474,67 \text{ DPD1} - 21.660,52 \text{ DPD2} - 14.753,86 \text{ DPD3} + 9.433,46 \text{ JART} + 571,19 \text{ JA5T} \\ - 14.967,72 \text{ DPK1} - 21.677,86 \text{ DPK2} - 18.340,26 \text{ DPK3} - 12.413,49 \text{ DPK4} - 9.292,64 \text{ DPK5}$$

3.2 Testing Assumptions

Several assumption tests are applied in Tobit regression analysis:

1. Normality Test

a. Hypotheses:

H_0 : The residuals are normally distributed.

H_1 : The residuals are not normally distributed.

b. Significance Level:

Alpha = 5%

c. Test statistics:

The Jarque-Bera (JB) test is used. The formula is:

$$JB = n \left(\frac{S^2}{6} + \frac{(K-3)^2}{24} \right) = 10.875$$

d. Rejection region:

Reject H_0 if $JB \geq \chi^2_{(\alpha; db=k)}$ or if $p\text{-value} < \alpha$

e. Conclusion:

Since $JB = 10.875 > \chi^2_{(\alpha; db=10)} = 18,307$ but $p\text{-value} = 2,2 \times 10^{-16} < \alpha = 0,05$ we reject H_0 . This means the residuals are **not normally distributed**.

2. Multicollinearity Detection

Multicollinearity occurs when two or more independent variables in a regression model are highly correlated, indicating a strong linear relationship between them. This condition can distort the estimation of regression coefficients and reduce the reliability of statistical inferences.

A good regression model should exhibit low or no multicollinearity. Detection is commonly done using the **Variance Inflation Factor (VIF)**.

Tabel 3. Multicollinearity Detection Results

Variable	Nilai VIF
DPD1	4,0787
DPD2	2,9285
DPD3	3,4841
JART	1,0376
JA5T	1,0166
DPK1	5,2038
DPK2	1,2762
DPK3	2,6496
DPK4	3,7892
DPK5	1,3633

Table 3 shows that the VIF value < 10 , meaning that all variables do not have multicollinearity.

3. Heteroscedasticity Detection

The heteroscedasticity hypothesis test is performed using the **Breusch-Pagan Test** [9] :

a. Hypotheses:

H_0 : There is no heteroscedasticity in the residuals of the regression model (homoscedasticity).

H_1 : There is heteroscedasticity in the residuals of the regression model.

b. Significance level:

$\alpha = 5\%$

c. Test statistics:

$$\phi = \frac{1}{2}(ESS) = 15,438$$

d. Rejection region:

Reject H_0 if $\phi_{hit} > \chi^2_{(\alpha; db=k)}$ or $p\text{-value} < \alpha$

e. Conclusion:

since $\phi = 15,438 < \chi^2_{(\alpha; db=10)} = 18,307$ or $p\text{-value} = 0,1502 > \alpha = 0,05$, we accept H_0 , his means there is no heteroscedasticity in the residuals of the regression model, indicating that the residual variance is constant.

3.3 Parameter Testing

The steps in testing regression parameters are:

1. Simultaneous Parameter Testing

To determine whether all the independent variables collectively influence the model, parameter testing is conducted using the G test.

a. Hypotheses:

$H_0: \beta_1 = \beta_2 = \dots = \beta_j = 0$ (all independent variables have no effect)

H_1 : At least one $\beta_j \neq 0, j = 1, 2, 3, \dots, k$ (Alpha significance value 5%)

b. Significance level:

$\alpha = 5\%$

c. Test statistic:

$$\begin{aligned} G &= -2[\ln L_R - \ln L_F] \\ &= -2[-12.922,31 - (-12.866,25)] \\ &= 112,12 \end{aligned}$$

d. Rejection region:

Reject H_0 if $G \geq \chi^2_{(\alpha; db=k)}$ or $p\text{-value} < \alpha$

e. Conclusion:

Since $G = 112,12 > \chi^2_{(0,05;10)} = 18,307$ or $p\text{-value} = 2,2 \times 10^{-16} < \alpha = 0.05$ we rejected H_0 . This indicates that **at least one independent variable significantly affects the model**.

2. Partial Parameter Testing

Partial parameter testing identifies which individual variables significantly influence the model. The **Wald test** is used for this purpose.

Tabel 4. Partial Parameter Testing Results

Variable	Wald	$\chi^2_{(\alpha; db)}$	P-Value	Keputusan
Intercept	6,0362	3,8415	$1,58 \times 10^{-09}$	*** H_0 rejected
DPD1	-4,3342	3,8415	$1,46 \times 10^{-05}$	*** H_0 rejected
DPD2	-2,9798	3,8415	0,00288	** H_0 rejected
DPD3	-2,2937	3,8415	0,02181	* H_0 rejected
JART	7,2699	3,8415	$3,60 \times 10^{-15}$	*** H_0 rejected
JA5T	0,1916	3,8415	0,84804	H_0 accepted
DPK1	-2,0060	3,8415	0,04485	* H_0 rejected
DPK2	-1,5015	3,8415	0,13322	H_0 accepted
DPK3	-2,2105	3,8415	0,02707	* H_0 rejected
DPK4	-1,6490	3,8415	0,09915	H_0 accepted
DPK5	-0,7389	3,8415	0,45995	H_0 accepted

Significance code: 0,001 = ***, 0,01 = **, 0,05 = *, 0,1 = .

Based on Table 4, the results of the Wald test indicate that six independent variables have a significant effect on the model at the 5% significance level. These variables are: Dummy Education 1 (household head graduated from elementary school), Dummy Education 2 (household head graduated from junior high school), Dummy Education 3 (household head graduated from high school), Number of household members, Dummy Job 1 (household head working in agriculture), Dummy Job 3 (household head working in trade). Conversely, four variables were found to be not statistically significant: Number of household members under 5 years old, Dummy Job 2 (household head working in mining), Dummy Job 4 (household head working in the service sector), Dummy Job 5 (household head working in the education sector).

3.4 Evaluation Of Model Goodness

Evaluation of the model's goodness is conducted by analyzing the error between the predicted and actual values as a percentage. This study employed the Mean Arctangent Absolute Percentage Error (MAAPE) to assess the forecasting model's accuracy. A lower MAAPE value indicates better predictive performance. In this case, the MAAPE value is 1.38%, which signifies that the prediction is highly accurate and represents a very reliable model outcome.

3.5 Interpretation

Based on the tobit regression assumption test, simultaneous test, and partial test, the result as follows:

- Intercep (51.246,32), This represents the average predicted household expenditure on egg and milk consumption when all independent variables are zero, assuming the data is not censored. It serves as the baseline expenditure value.
- Education Dummy Variable Coefficient, These variables reflect the influence of the household head's education level compared to those who did not complete elementary school:
 - DPD1 (-29.474,67) : household with heads who completed elementary school spend Rp 29,474.67 less, on average, than those without an elementary education. This suggests a decrease in both the likelihood of uncensored expenditure and the amount spent when uncensored.

- b. DPD2 (-21.660,52) : Junior high school graduates spend Rp 21,660.52 less on average. The negative effect is smaller than DPD1 but still significant.
 - c. DPD3 (-14.753,86) : High school graduates spend Rp 14,753.86 less on average. This also indicates a lower probability of being uncensored, though to a lesser extent.
3. Household Member Coefficient (JART +9.433,46) : each additional household member increases average consumption expenditure by Rp 9.433,46. This suggests a greater likelihood of uncensored expenditure and higher average spending among uncensored households.
4. Children Under 5 Years Coefficient (JA5T +571,19), Each additional child under five increases spending by Rp 571.19, but this effect is not statistically significant, indicating a limited impact on both censoring probability and average spending.
5. Occupational dummy coefficient, These reflect the impact of the household head's occupation compared to being unemployed:
 - a. DPK1 (-14.967,72), Agricultural workers spend significantly less, implying a higher chance of censored data or lower spending when uncensored.
 - b. DPK2 (-21.677,86), Mining workers also spend less, but the effect is not statistically significant.
 - c. DPK3 (-18.340,26), Those working in trade spend less significantly, suggesting a lower probability of expenditure above zero.
 - d. DPK4 (-12.413,49), Service sector workers spend less, though the effect is not significant.
 - e. DPK5 (-9.292,64), Education sector workers also show a lower, non-significant expenditure.

4 CONCLUSION

Based on the results and discussion of the Tobit regression analysis, the following conclusions can be drawn:

1. The Tobit regression model for household expenditure on egg and milk consumption in Bengkulu City is specified as follows:

$$y^* = 51.246,32 - 29.474,67 \text{ DPD1} - 21.660,52 \text{ DPD2} - 14.753,86 \text{ DPD3} + 9.433,46 \text{ JART} \\ + 571,19 \text{ JA5T} - 14.967,72 \text{ DPK1} - 21.677,86 \text{ DPK2} - 18.340,26 \text{ DPK3} \\ - 12.413,49 \text{ DPK4} - 9.292,64 \text{ DPK5}$$

2. The variables that significantly influence household expenditure on egg and milk consumption are : **Education Dummy 1** (elementary school graduate), **Education Dummy 2** (junior high school graduate), **Education Dummy 3** (high school graduate), Number of Household Members (JART), **Occupation Dummy 1** (household head working in agriculture **Occupation Dummy 3** (household head working in trade)

REFERENCES

- [1] Hosmer, D. W. dan S Lemeshow. 2000. Applied Logistic Regression Second Edition. John Wiley & Sons, Inc, New York.
- [2] M. A. Alwansyah, "Survival Analysis of Students Not Graduated on Time Using Cox Proportional Hazard Regression Method and Random Survival Forest Method," J. Stat. Data Sci., vol. 2, no. 1, pp. 13–21, 2023.
- [3] Tobin, J. 1958. Estimation of Relationship for Limited Dependent Variabels. *Econometrica*, vol 26.
- [4] Sinurat, E., S. Nugroho., E. Sunandi. 2014. Analisis Regresi Tobit (Studi kasus: Faktor-Faktor Yang Berpengaruh Terhadap Biaya Pengeluaran Konsumsi Pakaian Dikalangan Mahasiswa Matematika Angkatan 2010-2013 FMIPA Unib). *Jurnal Jurusan Matematika FMIPA UNIB*.
- [5] Suhardi, I. Y. dan R. Llewelyn. 2001. Penggunaan Model Regresi Tobit untuk Menganalisa Faktor-faktor yang Berpengaruh Terhadap Kepuasan Konsumen untuk Jasa Pengangkutan Barang. *Jurnal Manajemen & Kewirausahaan*. 3.
- [6] Gujarati, D. 2004. Basic Econometrics, Fourth Edition. Mc-Graw-Hill, Inc, New York.
- [7] Akinleye, O., Agboola, J., Isaac, A., dan Oluwaseun, A.-I. 2020. Comparison Of Different Tests For Detecting Heteroscedasticity. December 2020.
- [8] Kim, S., & Kim, H. 2016. A new metric of absolute percentage error for intermittent demand forecasts. *International Journal of Forecasting*, 32(3), 669–679.
- [9] Breusch, & Pagan. 1979. A Simple Test for Heteroscedasticity and Random Coefficient Variation. *Econometrica*, 47(5), 1287–1294