

Generalized Linear Model analysis for Binomial distribution with the application

Dr. Kamran Hassan Ahmed

Department of Statistics and Informatics, College of Administration and Economics,
University of Salahaddin-Erbil, Kurdistan Region, Iraq.

Email: kamarn.ahmed@su.edu.krd

ARTICLE INFO

Article History:

Received: 16/8/2023

Accepted: 10/10/2023

Published: Winter 2024

Keywords:

Generalized Linear

Model, Binomial

Logistic Regression,

Probit & Logit Models.

Doi:

10.25212/lfu.qzj.9.4.47

ABSTRACT

The aim of this research is to use a generalized Linear model for the binary logistic regression model to study the effect of different concentrations of two drugs (x_{i1} , x_{i2}) that were studied at different levels on patients with nephritis (kidney inflammation), where the response variable (y_i) represents the number of cure cases (Binomial distribution) as a result of taking the two drugs. In order to estimate the logistic model, two functions link were used, the first is the probit function, and the second is the logit function, and then a comparison between the results of using the two functions. The results support the preference of the estimated Logit regression model through the criteria (Akaike's Information Criterion, Bayesian Information Criterion, and Mean Square Error). The effectiveness of the first independent variable (drug) on the second independent variable (drug) through the preferred model.

1. Introduction

Regression analysis is a model that analyzes and explains the relationships between the response variable and the independent variables by relating these variables to a mathematical equation that may be linear or nonlinear. After determining the form of this relationship, the model's parameters are estimated for interpretation or prediction, according to the nature of the study.) Logistic regression is a type of regression in which the response variable is a qualitative variable that may take two

values (binary logistic regression) and may take more than two values (multiple logistic regression) (Ali, 2011). In logistic regression, our goal is not to explain the change in the values of the response variable (Ali & Tara 2010), but the interpret the probability of occurrence and non-occurrence of the phenomenon under study, and represent the logistic regression equation. In this research, binary logistic regression will be used to study the effect of different concentrations of two studied drugs at different levels on patients with inflammation of the kidneys. To estimate the logistic model, two link functions were used (probit and logit function), and then their results were compared.

2. Theoretical Aspect

2.1. Generalized Linear Model

Generalized linear models (GLM) are extensions of linear models in which the response variables are linked to the linear model using a link function (Ali et al 2023). Additionally, the model permits a non-normal distribution for the response variable. It covers popular statistical models like logistic models for binary data, linear regression for normally distributed responses, and linear models for counting data (such as binomial distribution, Poisson distribution, and others) (Menard & Scott, 2002).

Generalized linear models are linear models where the response variable is modelled by a linear function of the independent variables (Omar et al 2020). GLMs consist of the following three elements:

1. The random element: The probability distribution of the response variable is the random element in a GLM.
2. The systematic element: includes the independent variables that linearly combine to the predictor of the GLMs.
3. The link function: The random and systematic components of the GLMs are linked or related by the link function. It explains how the values of the response variables relate to the independent variables of the linear predictor (Ali & Saleh 2022).

2.2. Binomial Logistic Regression

The logistic model is used when the response variable has a binomial distribution through the general linear model that does not assume the normal distribution of random error so that the response variable is linear combination to the independent variables through a specific link function such as Logit (Shahla, 2023).

For a binomial distribution, we assume that there are frequencies for each level of the independent variable (or several p-level independent variables), meaning that they have levels (x_1, x_2, \dots, x_m) and that there are n_i values of (y) observations that may be repeated for each level of x where (y) takes values zero or one (Raza et al 2018), n_i times for each level of (x) , if we assume that (y_i) is the number of appearances of one for each level of (x) , then the percentage of appearance of one in each level of x is (Ali & Tara 2007).

$$p_i = \frac{y_i}{n_i}, \quad i = 1, 2, \dots, m \quad (1)$$

On this basis, the logistic function can be obtained by using the transformation on the observed ratios, i.e. (Hamad, 2016):

$$p_i^* = Ln\left(\frac{p_i}{1 - p_i}\right) \quad (2)$$

Since the random error variance is not homogeneous, so the weighted least squares method can be used to estimate the parameters of the logistic model, and the weight in this case is the inverse of the variance, meaning that (Omar et al 2020):

$$\hat{w}_i = \frac{1}{v(p_i^*)} \quad (3)$$

The ratio in the sample (\hat{p}_i) represents a parameter of the binomial distribution with the mean (p_i) and the variance $p_i(1 - p_i) / n_i$, where (p_i) is the probability that the studied trait affected by the independent variable has a certain level (x_i) , and that the relationship between (p_i) (Which represents the qualitative response variable or dummy) and (x_i) which represents the independent variable in the logistic model depends on the following function (Agresti, 1990):

$$p_i = f(x_i) \quad (4)$$

The converted variable (p_i^*) has a mean equal to $\ln(p_i / (1 - p_i))$ and the variance is the reciprocal of the variance of a binomial distribution, meaning that:

$$V(p_i^*) = \frac{1}{n_i p_i (1 - p_i)} \quad (5)$$

Therefore, the weights that can be used in the weighted least squares method can be estimated as follows (Bradley, 1997):

$$\hat{w}_i = n_i p_i (1 - p_i); \quad i = 1, 2, \dots, m \quad (6)$$

On this basis, the regression parameters can be estimated using the weighted least squares method and obtaining the linear logistic function (Hosmer & Lemeshow, 2000):

$$\hat{p}_i^* = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip} \quad (7)$$

Through this equation, the estimated values (\hat{p}_i^*) can be found by substituting for the levels of (x_i) values for all independent variables, and then converting the (\hat{p}_i^*) values to their original value through the following formula:

$$\hat{p}_i = \frac{\text{Exp}(\hat{p}_i^*)}{1 + \text{Exp}(\hat{p}_i^*)} \quad (8)$$

2.3. Probit & Logit Models

The probit and logit models are used to model dichotomous or binary response variables in statistical modeling. If our outcome is dichotomous (Ali, 2022), the natural distribution to consider for a *GLM* is the binomial, $y \sim \text{Binomial}(np)$ with (p) being the mean of the binomial, and (n) being the number of trials. Logit models, also known as logistic regression models, are a category of statistical models that are used to

estimate the likelihood that an event will occur (Ali & Jwana, 2022). It is used to model an event's odds of success as a function of explanatory variables. Mathematically it can be written as below

$$\text{logit}(I) = \log\left(\frac{p}{1-p}\right) = Z = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} \quad (9)$$

where p is the probability that an event occurring, and one is the odds that it will. From this, we may also calculate the likelihood that the events will occur.

$$p = \pi(Z) = \left(\frac{1}{1 + \exp(-Z)}\right) \quad (10)$$

Where $\pi(Z)$ is a logistic function. As the value of Z approaches $-\infty$ the value of $\pi(Z)$ or p approaches 0. Additionally, when the value of Z becomes closer to $+\infty$, the value of $\pi(Z)$ or p gets closer to 1.

Probit model is similar to logit model, it determines the possibility that an item or event will fall into one of a number of categories by calculating the probability that an observation with a certain attribute will fall into a specified category. The following formula can be used to represent the Probit model:

$$P(y = 1|x) = \Phi(Z) = Z = \Phi(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \quad (11)$$

Where y is the likelihood that the event will happen (hence, $y = 1$). Φ is the cumulative standard normal distribution function. Z is linearly related of explanatory variables (x). Logistic function is used instead of the cumulative standard normal distribution function (Φ) in the case of the logit model. The figure below represents the Probit & Logit models:

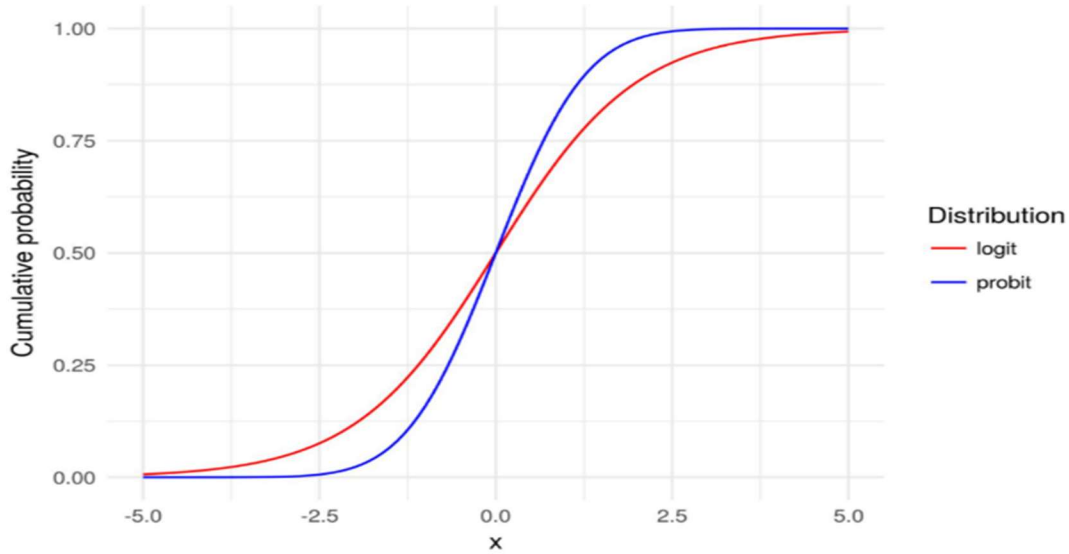


Figure 1. Predictive values for Probit and Logit Regression Models

2.4. Evaluation criteria

Akaike's Information Criterion (*AIC*), Bayesian Information Criterion (*BIC*), and Mean Square Error (*MSE*) will be used (Kareem, 2020) as evaluation criteria to evaluate the estimators of binary logistic regression models.

$$AIC = -2(\text{LogLikelihood}) + 2k \quad (12)$$

$$BIC = -2(\text{LogLikelihood}) + 2k\text{Log}(n) \quad (13)$$

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{(n - 1)} \quad (14)$$

Where *k* is number of estimated parameters and *n* is sample size, and the lowest value is for the best model (Ali & Saleh 2022).

3. Application Aspect

The effect of different concentrations of two drugs (two independent variables) at different levels (Milligrams per dose) x_{r1} , x_{r2} was studied on patients with nephritis (inflammation of the kidneys), the response variable (Y_i) represents the number of

cases of recovery as a result of taking the two drugs through (15) observations representing ($N = 15$) samples (n_i) of different size taken from one the hospitals were as follows:

Table 1. Data of Nephritis patients at different levels (Milligrams per dose)

Samples	Y_i	x_{i1}	x_{i2}	n_i
1	23	.51	.39	60
2	21	.55	.28	50
3	80	.67	.18	100
4	60	.63	.15	80
5	20	.23	.54	100
6	25	.33	.44	110
7	30	.44	.35	115
8	35	.60	.33	95
9	42	.83	.22	80
10	45	.62	.42	65
11	29	.70	.69	40
12	14	.75	.55	45
13	18	.82	.67	30
14	34	.77	.72	70
15	30	.78	.66	48

The response variable (categorical) information is summarized in Table .2:

Table 2. Categorical Variable Information

		n	Percent
Response Variable	Y	Events	506 46.5%
		Non-Events	582 53.5%
		Total	1088 100.0%

The categorical (response) variable information table shows that there were (506) cases of cure, with a rate of 46.5%, compared to (582) cases of non-recovery, with a rate of 53.5%, for the total samples examined, which amounted to (1088) cases. The information on the independent variables is summarized in Table .3:

Table 3. Continuous Variable Information

		N	Minimum	Maximum	Mean	Std. Deviation
Covariate	x1	15	.23	.83	.6153	.17804
	x2	15	.15	.72	.4393	.19144

The homogeneity of the random error variance is first tested and summarized in Table 4.

H_0 : The random error variance homogeneity

H_1 : The random error variance heterogeneity

Table 4. Tests of Homogeneity of Variances

Levene Statistic	Df1	Df2	Sig.
5.084	1	13	0.042

The homogeneity test supports the hypothesis of random error variance heterogeneity (Mustafa & Ali, 2013), which is consistent with the theoretical aspect. Therefore, weights that depend on formula (5) were used, and then the homogeneity test was re-tested as in Table .5:

Table 5. Tests of Homogeneity of Variances

Levene Statistic	Df1	Df2	Sig.
1.537	1	13	0.237

The homogeneity test supports the random error variance homogeneity hypothesis, so the analysis was carried out.

3.1. Probit Regression Model

The goodness of fit Logit regression model is summarized in Table .6:

Table 6. Goodness of Fit (Probit)

Criteria	Value
Akaike's Information Criterion (AIC)	2367.151
Bayesian Information Criterion (BIC)	2382.127
Mean Square Error (MSE)	128.3975

The lower the criteria AIC, BIC, and MSE, the better the estimated model. The Omnibus test table tests the significance of the model parameters as a whole except for parameter (β_0), i.e., testing the following hypothesis:

$$H_0 : \beta_i = 0 \quad \text{for all } i$$

$$Vs \quad H_1 : \beta_i \neq 0 \quad \text{for at least 1 coefficient } t$$

Table 7. Omnibus Test (Probit)

Likelihood Ratio Chi-Square	df	Sig.
1776.317	2	.000

Table. 7 show that the value χ^2 is equal to (1776.317) for the model, which is greater than its tabular value under the level of significance (0.01) and degrees of freedom (2) which equals (10.60). That is, rejecting the null hypothesis and accepting the alternative hypothesis, which means that there is at least one of the two parameters that are not equal to zero, that is, the existence of a variable with at least one independent has a significant effect on the response variable and this is confirmed by the p-value which is equal to zero and is less than the level of significance (0.01). The significance test of the estimated model parameters is summarized in Table. 8:

H_0 : The model does not fit the data

H_1 : The model is fit to the data

Table 8. Tests of Model Effects (Probit)

	Wald Chi-Square	df	Sig.
(Intercept)	551.331	1	.000
x1	1367.872	1	.000
x2	223.403	1	.000

The table Tests of model effects (Table. 8) show that the values of Wald (chi-square) were greater than the value of chi-squared under the significance level (0.01) and degrees of freedom (1) equal to (7.88), which indicates the significance of the parameters estimated for each of them with the stability of the effect of the rest, and

this is confirmed by the p-values, which are equal to zero and are less than the significance level (0.01). The estimated parameters table is shown in Table. 8:

Table. 9. Parameter Estimates (Probit)

Parameter	B	Std. Error	99% Wald Confidence Interval		Exp(B)
			Lower	Upper	
(Intercept)	-.975	.0415	-1.082	-.868	.377
x1	2.053	.0555	1.910	2.196	7.795
x2	-.878	.0588	-1.030	-.727	.416

The table of estimated parameters shows that the logistic regression model for the binomial distribution is as follows:

$$\text{Probit} [P (x)] = -0.975 + 2.053x_{i1} - 0.878x_{i2}$$

The values of the coefficients on the likelihood scale Exp(B) mean that each one-unit increase in the independent variable (the concentration of the first drug, for example) is associated with an increase in the likelihood of the response variable (recovery from disease) by (7.795) with the neutralization of other variables.

3.2. Logit Regression Model

The goodness of fit Logit regression model is summarized in Table. 10:

Table 10. Goodness of Fit (Logit)

Criteria	Value
Akaike's Information Criterion (AIC)	2357.199
Bayesian Information Criterion (BIC)	2372.175
Mean Square Error (MSE)	126.8175

The lower the criteria AIC, BIC, and MSE, the better the estimated model.

Table 11. Omnibus Test (Logit)

Likelihood Ratio Chi-Square	df	Sig.
1786.269	2	.000

The results of the Omnibus test Table. 11 support the goodness of fit of the data for the estimated model. The significance test of the estimated model parameters is summarized in Table (12):

Table 12. Tests of Model Effects (Logit)

	Wald Chi-Square	df	Sig.
(Intercept)	556.317	1	.000
x1	1296.695	1	.000
x2	234.152	1	.000

The table tests of model effects (Table. 12) show that the regression model coefficients are all significant. The estimated parameters table is shown in Table. 13:

Table 13. Parameter Estimates

Parameter	B	Std. Error	99% Wald Confidence Interval		Exp(B)
			Lower	Upper	
(Intercept)	-1.605	.0680	-1.780	-1.430	.201
x1	3.387	.0941	3.145	3.630	29.587
x2	-1.452	.0949	-1.696	-1.207	.234

The table of estimated parameters shows that the logistic regression model for the binomial distribution is as follows:

$$Ln(\hat{Odds}) = -1.605 + 3.387x_{i1} - 1.452x_{i2}$$

The values of the coefficients on the likelihood scale Exp(B) mean that each one-unit increase in the independent variable (the concentration of the first drug, for example) is associated with an increase in the likelihood of the response variable (recovery from disease) by (29.587) with the neutralization of other variables.

3.3. Comparison of Probit and Logit regression models

For comparison between the two models, the efficiency criteria of the estimated models are summarized in Table. 14.

Table 14. Comparison of Probit and Logit regression models criteria

Method	MSE	AIC	BIC	Likelihood Ratio Chi-Square
Probit	128.3975	2367.151	2382.127	1776.317
Logit	126.8175	2357.199	2372.175	1786.269

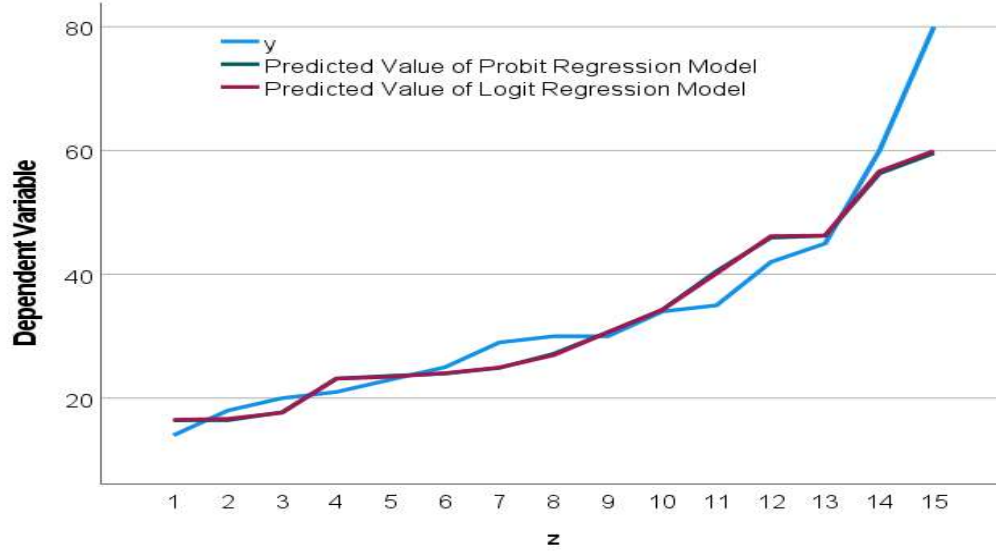


Figure .2. Predictive values for Probit and Logit Regression Models

The results of the criteria comparing Table. 13 of the Probit and Logit regression models support the preference of the estimated Logit regression model because the values of the criteria for MSE, AIC, and BIC were less, while the value of Likelihood Ratio (Chi-Square) was greater (better). Therefore, the estimated Logit regression model is relied upon. The increase in the likelihood of the response variable (recovery from disease) by (29.587) for the first independent (drug) variable (x_1). While there is a very small increase in the number of healing cases due to the independent (drug) variable (x_2) and its amount (0.234).

Figure (2) shows the real and predictive values (for the response variable) for the number of healing cases at standard levels (z) from (1-15). Where the figure shows

that there is a great convergence between the estimates of the Probit and Logit regression models, as shown on the theoretical side.

4. Conclusions and recommendations

4.1 Conclusions

The following is a presentation of the most important conclusions based on the results:

- 1- Significance of both models (Probit and Logit) according to the Likelihood Ratio Chi-Square test
- 2- The estimated logit regression model is better than the probit regression model because the values of MSE, AIC and BIC criteria were lower, while the probability ratio value (Chi-Square) was greater (better). Therefore, the estimated logit regression model is more reliable.
- 3- Both medications have a significant effect on increasing the rate of recovery from nephritis
- 4- The first medication has a greater effect than the second medication in increasing the recovery rate.

4.2 Recommendations

The following is a presentation of the most important recommendations:

- 1- There is a convergence between the estimates of the probit and logit regression models, but our recommendation is to use the logit regression models.
- 2- The use of the logistic regression model in other fields (social and economic) and not only in the medical fields
- 3- Conducting other studies by introducing new variables and factors affecting kidney disease.

References:

- 1- Agresti, A., (1990): "Categorical Data Analysis" John Wiley Sons. Inc. New York.
- 2- Ali, Taha Hussein and Jwana Rostam Qadir. "Using Wavelet Shrinkage in the Cox Proportional Hazards Regression model (simulation study)", *Iraqi Journal of Statistical Sciences*, 19, 1, 2022, 17-29.
- 3- Ali, Taha Hussein, "Estimation of Multiple Logistic Model by Using Empirical Bayes Weights and Comparing it with the Classical Method with Application" *Iraqi Journal of Statistical Sciences* 20 (2011): 348-331.
- 4- Ali, Taha Hussein, and Dlshad Mahmood Saleh. "COMPARISON BETWEEN WAVELET BAYESIAN AND BAYESIAN ESTIMATORS TO REMEDY CONTAMINATION IN LINEAR REGRESSION MODEL" *PalArch's Journal of Archaeology of Egypt/Egyptology* 18.10 (2021): 3388-3409.
- 5- Ali, Taha Hussein, and Saleh, Dlshad Mahmood, "Proposed Hybrid Method for Wavelet Shrinkage with Robust Multiple Linear Regression Model: With Simulation Study" *QALAAI ZANIST JOURNAL* 7.1 (2022): 920-937.
- 6- Ali, Taha Hussein. "Modification of the adaptive Nadaraya-Watson kernel method for nonparametric regression (simulation study)." *Communications in Statistics-Simulation and Computation* 51.2 (2022): 391-403.
- 7- Ali, Taha Hussein; Tara Ahmed Hassan. "A comparison of methods for estimating regression parameters when there is a heterogeneity problem of variance with a practical application", *Journal of Economics and Administrative Sciences*, 16.60 (2010): 216-227.
- 8- Ali, Taha Hussein; Tara Ahmed Hassan. "Estimating of Logistic Model by using Sequential Bayes Weights", *Journal of Economics and Administrative Sciences*, 13.46 (2007): 217-235.
- 9- Bradley, Andrew P. (1997): "The Use of the Area under the ROC Curve in the Evaluation of Machine Learning Algorithms. *Pattern Recognition*" Vol.30, No.7, pp.1145-1159.
- 10- Hamad, Karzan Faidh. (2016): "Applying Logistic Regression Model to Study Some Statistical Measurements in Distinguishing Shapes of Mass in Medical Images" B.Sc. *Statistical Sciences- Salahaddin University*.
- 11- Hosmer, David W. & Lemeshow, Stanley (2000): "Applied Logistic Regression" 2nd Edition, Johnson Wiley & Sons Incorporation, New York, USA.
- 12- Kareem, Nazeera Sedeek, Taha Hussein Ali, and Awaz Shahab M, "De-noise data by using Multivariate Wavelets in the Path analysis with application", *Kirkuk University Journal of Administrative and Economic Sciences*, 10.1 (2020): 268-294.

- 13- Menard, Scott. (2002): *"Applied Logistic Regression Analysis" Second Edition, Sage Publication, Inc.*
- 14- Mustafa, Qais, and Ali, Taha Hussein. *"Comparing the Box Jenkins models before and after the wavelet filtering in terms of reducing the orders with application." Journal of Concrete and Applicable Mathematics 11 (2013): 190-198.*
- 15- Omar, Cheman, Taha Hussien Ali, and Kameran Hassan, *Using Bayes weights to remedy the heterogeneity problem of random error variance in linear models, IRAQI JOURNAL OF STATISTICAL SCIENCES, 17, 2, 2020, 58-67.*
- 16- Raza, Mahdi Saber, Taha Hussein Ali, and Tara Ahmed Hassan. *"Using Mixed Distribution for Gamma and Exponential to Estimate of Survival Function (Brain Stroke)." Polytechnic Journal 8.1 (2018).*
- 17- Shahla Hani Ali, Heyam A.A.Hayawi, Nazeera Sedeek K., and Taha Hussein Ali, (2023) *"Predicting the Consumer price index and inflation average for the Kurdistan Region of Iraq using a dynamic model of neural networks with time series", The 7th International Conference of Union if Arab Statistician-Cairo, Egypt 8-9/3/2023:137-147.*

شيكارى مۆدېلى ھېلى گشتگير بۆ دابه شېبوى دووانه يى له گه ل به كارھېنان

پوخته:

ئامانجى ئەم توپزېنە وە يە بە كارھېنانى مۆدېلىكى ھېلى گشتگيرە بۆ مۆدېلى لار بونە وە يى لوجستىكى دووانه يى بۆ ليكۆلېنە وە له كارىگەرى دوو جۆرى دەرمان له چەند ئاستىكدا لەسەر نەخۆشانى تووشبوو بە ھەوکردنى گورچيلە، كە گۆراوى وابەستە برىتېيە له ژمارەى حالەتەكانى چارەسەرکردن (دابه شېبوى دووانه يى) له ئەنجامى بە كارھېنانى دوو دەرمانەكە. بۆ خەملاندنى مۆدېلى لوجستى، دوو نەخشەى بەستەنە وە بە كارھېنان، يەكەمىان برىتېي بوو له نەخشەى (Probit) دەكات، دووھەمىان برىتېي بوو له نەخشەى (Logit) دەكات، پاشان بەراوردکردنى ئەنجامى بە كارھېنانى دوو نەخشەكە. ئەنجامەكان پشتگيرى له كارايى مۆدېلى لار بونە وە يى (Logit) دەكەن له رېگەى پېوھەرەكانى بە كارھېنراو (پېوھەرەكانى زانيارى ئاكاكى، پېوھەرەكانى زانيارى بېز و مامناوھەندى ھەلەى چوارگۆشە) ھەروھە پەسەندکردنى دەرمانى يەكەم له چارەسەرکردنى نەخۆشانى ھەوکردنى گورچيلە بە بەراورد بە دەرمانى دووھەم بە پېى مۆدېلى خەملاينراو.

تحليل الأنموذج الخطي المعمم لتوزيع ذي الحدين مع التطبيق

الملخص:

يهدف البحث الى استخدام الأنموذج الخطي المعمم لنموذج الانحدار اللوجستي الثنائي في دراسة تأثير نوعان من الأدوية عند عدة مستويات على مرضى الالتهاب الكلوي، حيث المتغير التابع يمثل عدد حالات الشفاء (توزيع ذي الحدين) نتيجة استخدام الدوائيين. لتقدير الأنموذج اللوجستي تم استخدام دالتي الربط تمثل الأولى دالة (Probit) في حين تمثل الثانية دالة (Logit) ومن ثم المقارنة بين نتائج استخدام الدالتين. تدعم النتائج كفاءة نموذج الانحدار (Logit) من خلال المعايير المستخدمة (معيار معلومات أكايكي، معيار معلومات بيز ومتوسط الخطأ التريبيعي) فضلاً عن أفضلية الدواء الأول في علاج مرضى الالتهاب الكلوي مقارنة بالدواء الثاني من خلال أداء الأنموذج المقدر.