

# Application of Logistic Regression on Heart Disease Data and a Review of Some Standardization Methods

Florence George<sup>1</sup>, Sultana Mubarika Rahman Chowdhury<sup>1</sup> & Sneh Gulati<sup>1</sup>

<sup>1</sup> Mathematics and Statistics Department, Florida International University, Miami, Florida, USA

Correspondence: Florence George, Mathematics and Statistics Department, Florida International University, Miami, Florida, 33199, USA

Received: July 17, 2022 Accepted: August 25, 2022 Online Published: August 29, 2022

doi:10.5539/ijsp.v11n5p1

URL: <https://doi.org/10.5539/ijsp.v11n5p1>

## Abstract

The purpose of this study is to do a review of logistic regression and its applications. In addition to the review, a comparison of four different methods of standardization of the  $\beta$  - coefficients was done using publicly available Heart Disease Data. The methods were compared using their performance in testing accuracy, training accuracy, and area under the curve (AUC). Based on the comparisons, it was evident that standardizing the coefficient did not affect the overall prediction accuracy of the model regardless of the method used. Although there was some difference found in the training and testing accuracies, the AUC's were similar to the unstandardized model for all methods. In essence, standardizing facilitates better interpretation and does not affect the predictive accuracy of the model.

**Keywords:** Logistic Regression, Logit model, Standardized Coefficients

## 1. Introduction

Logistic regression analysis is a specialized case of regression analysis, where the variable to be predicted is classified into two or more categories. In such cases, the traditional regression technique fails to explain the association between the independent variables and the response variable. Binary logistic regression model or logit model is the most common form of this method of analysis in which the response variable takes only two values (Menard, 2000).

The specific form of a binary logistic regression model generally used is

$$P(Y = 1) = \frac{e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p)}}, \quad (1)$$

where  $Y$  is the dependent variable and  $X_1, X_2, \dots, X_p$  are the independent variables. The dependent variable  $Y$  takes on the values either 0 or 1; where 1 indicates the occurrence of a specific event and 0 indicates the absence. Therefore,  $P(Y = 1)$  represents the probability of that event happening and  $P(Y = 0)$  depicts the probability of the event is absent.

Logistic regression has a wide range of applications in various fields and its functionality has increased dramatically in the past several decades. While multiple linear regression falls short in analyzing data with response variable that is not continuous, logistic regression gives an essential tool in such cases. Application of this method is not limited to only binary cases as it can be easily modified for cases where response variables have more than two categories. Risk factor analysis and predictive modeling is one of the main implementations of logistic regression (Peterson, L. E. et al., 1995). Logistic regression can also be used in survival analysis by grouping event times into intervals and converting them to categories (Abbott, 1985). Hence, is broadly used in medical research fields to examine the association between risk factors and diseases (Kurt, I. et al., 2008; Hassanipour, S. et al., 2019).

The parameters, the standard error of the parameters, and the measures of the goodness of fit are estimated using the methods of maximum likelihood estimation (Greene, 1993, Peng et al., 2002)

The logit transformation of  $P(Y = 1)$  is defined as

$$\text{logit}(Y) = \ln\left(\frac{P(Y = 1)}{P(Y = 0)}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p. \quad (2)$$

The model converts the nonlinear relationship between  $P(Y = 1)$  and the independent variables to a linear equation that explains the effect they may have on the dependent variable. This linear form gives the opportunity to interpret the coefficients of the proposed model.

The interpretation of results is rendered using the widely used odds ratio technique for both categorical and continuous predictors (Peng et al., 2002). Even though the odds ratio can give an idea of the direction of the relationship between the response variable and explanatory variables, it is not enough to explain the overall extent of how they are related and also it falls short of comparing over models (Allison, 1999). It should also be noted that some alternate methods based on the effect measures are proposed in several papers to explain the effects of covariates on binary response variables in a logistic regression model (Agresti A., Kateri M., 2017; Agresti A., Tarantola C., 2018).

However, the primary focus of this paper is on the  $\beta$ -coefficients and does not investigate these alternate methods. Standardizing the  $\beta$ -coefficients is another approach found in various literature studies (Long J. S., 1997; Menard S.W., 1995), and different techniques to standardize the  $\beta$ -coefficients have also been proposed to allow for more meaningful interpretations. Standardized coefficients become invariant to the change in scale of measurement which enables one to compare the relative influence of different explanatory variables within logistic regression (Agresti A., 2018; Agresti A., Finlay B., 1997). However, even though there are some proposed standardized, semi-standardized coefficients for logistic regression none of them can be universally defined. Robert L. Kaufman (Kaufman R.L., 1996) in his study found that semi-standardized coefficients measuring the change in predictive probability of outcomes are preferable because they are intuitively appealing and as they are bounded in the interval  $[-1, +1]$ , interpretation of their magnitude becomes easier. Some approaches of standardizing the coefficients were analyzed using a practical example by Scott Menard, which included both semi-standardized and completely standardized techniques (Menard S., 2004).

In this paper, we will discuss the four methods discussed by Menard and in addition to that, we propose a modification of these four methods for standardization of logistic regression coefficient. These methods will also be compared based on the resulting testing accuracy, training accuracy, AUC (area under the curve).

The simplest method of partial standardization of logistic regression coefficients is to multiply the coefficients by their individual standard deviation. This method was mentioned by Menard (Menard S.W., 1995).

$$b_1 = b * S_x, \quad (3)$$

where, the standard deviation of the explanatory variable  $X$  ( $S_x$ ) is multiplied with the unstandardized estimated coefficient of the corresponding variable  $b$ . This can be considered as the only predictor-based standardization technique. Another similar approach is to change the scale of both the dependent variable and the predictors using the standard deviation of the standard logistic distribution. That is,

$$b_2 = \frac{b * S_x}{\frac{\pi}{\sqrt{3}}}, \quad (4)$$

where,  $(\pi \sqrt{3}) = 1.8$ . This method has been adapted in SAS to standardize the coefficients in the PROC LOGISTIC procedure. Long suggested another approach for standardization which includes the standard deviation of the standard normal distribution (Long J. S., 1997).

The calculation of this method is similar to the previous one, the only difference is the standard deviation of the standard normal distribution is added with the standard deviation of the logistic distribution. Hence Equation (4) becomes,

$$b_3 = \frac{b * S_x}{\frac{\pi}{\sqrt{3}} + 1}. \quad (5)$$

All of these standardized coefficients only take into account the variation of the independent variable. Hence, they cannot be considered as fully standardized. To standardize the response variable standard deviation of *logit* ( $y$ ) needs to be calculated, which is tricky. A way out of this is to use the standardization followed in OLS, which is defined as follows,

$$b^{**} = b * \frac{S_x}{S_y}.$$

Again, from the definition of Coefficient of Determination ( $R^2$ ), we get

$$R^2 = \frac{S_{\hat{y}}^2}{S_y^2},$$

where,  $\hat{y}$  is the estimated value of  $y$ . Adjusting the equation for OLS we get,

$$S_y^2 = \frac{S_{\hat{y}}^2}{R^2},$$

Substituting  $\text{logit}(y)$  in case of  $y$  and  $\text{logit}(\hat{y})$  in the place of  $\hat{y}$  we get for logistic regression,

$$S_{\text{logit}(y)}^2 = \frac{S_{\text{logit}(\hat{y})}^2}{R^2}.$$

Hence, using the similar strategy used in OLS the estimated coefficients can be standardized as follows

$$b_4 = \frac{(b * S_x)(R)}{S_{\text{logit}(\hat{y})}}. \quad (6)$$

This coefficient can be considered as fully standardized as it also takes into account the variance of the response variable in contrast to the other coefficients discussed before where only the variation of the predictor was studied. For the purpose of comparing the above four standardization methods, they will be applied to z- scaled data using the mean and standard deviation. Since median and MAD may be better measures for scaling asymmetric data, we propose applying these standardization techniques to the median and the MAD scaled data.

In the next section, these standardized logistic regression coefficients for both z-scaled and median/MAD scaled data will be compared by applying the methods to Heart Disease Data.

## 2. Implementation of Standardization Methods

In order to illustrate the calculation of the standardization techniques and to review the outcomes, the Cleveland Heart disease dataset was used. It is a widely used dataset that is publicly available online (Detrano R., 1989). The aim was to apply logistic regression to develop a predictive model for heart diseases using the predictors. The four different coefficient standardization methods were applied to the coefficients of the customary model. After that, the resultant models were compared based on their prediction accuracy.

### 2.1 Dataset Details

Originally, the data set contained 76 attributes, but a subset of 14 variables are generally used by the researchers in all published experiments with a total of 313 observations. The 14 variables include a response variable "target" which refers to the presence of heart disease in the patient. For the target variable, a value of 0 indicates no/ less chance of heart attack while a value of 1 indicates yes/ high chance of heart attack.

The 13 predictors considered in the dataset are as follows (Detrano R., 1989):

1. *AGE*: Continuous
2. *SEX*: Categorical ( 0 = Female, 1 = Male)
3. *Chest Pain Type(CP)*: Categorical (4 values) 0: typical angina 1: atypical angina 2: non - anginal pain 3: asymptomatic
4. *Trestbps*: Continuous, represents resting blood pressure on admission
5. *Chol*: Continuous, represents Serum cholesterol in mg/dl
6. *Fbs*: Categorical , represents fasting blood sugar level, (2 values) 1: True - fasting blood sugar is greater than 120 mg/dl 0: False - fasting blood sugar is less than 120 mg/dl
7. *Restecg*: Categorical,represents resting electrocardiographic outcomes (4 values) 0: normal 1: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of >0.05 mV) 2: showing probable or definite left ventricular hypertrophy by Estes' criteria)
8. *Thalach*: Continuous, represents maximum heart rate achieved
9. *Exang*: Categorical, represents the existence of exercise-induced angina (2 values Yes/No)
10. *Oldpeak*: Continuous, ST depression induced by exercise relative to rest
11. *Slope*: Categorical, represents the slope characteristics of the peak exercise ST segment

12. *Ca*: Discrete, represents the number of fluoroscopy colored major vessels (values 0-3);

13. *Thal*: Categorical, (3 values) 0: normal 1: fixed defect 2: reversible defect

## 2.2 Methodology and Results

Primarily, logistic regression was applied to the complete dataset. Four standardization techniques of the coefficients discussed in the previous section were applied to this result. Calculation of  $b_1$  is done by simply multiplying the standard deviation of each explanatory variable with their corresponding coefficients. For instance, for *Age*  $b_1 = (-0.004908) * (9.0821010) = -0.04457922$  and so on Table ??.

Table 1. Modified coefficients using different standardization methods

	Customary model (Pvalue)	Method 1	Method 2	Method 3	Method 4
<b>Intercept</b>	3.4505				
<b>Age</b>	- 0.0049(0.8323)	- 0.0446	- 0.0246	- 0.0158	- 0.0110
<b>Sex</b>	- 1.7582 (0.0002)	- 0.8193	- 0.4517	-0.2912	- 0.2019
<b>Cp</b>	0.8599 (0.000)	-0.8874	0.4893	0.3154	0.2189
<b>Trestbps</b>	-0.0195 (0.0596)	-0.3416	-0.1883	-0.1214	-0.0842
<b>Chol</b>	-0.0046 (0.2209)	-0.2400	-0.1323	-0.0853	-0.0591
<b>Fbs</b>	0.0349 (0.9475)	0.0124	0.0069	0.0044	0.0031
<b>Restecg</b>	0.4663 (0.1806)	0.2452	0.1352	0.0871	0.06043
<b>Thalach</b>	0.0232 (0.0265)	0.5317	0.2931	0.1889	0.1310
<b>Exang</b>	-0.9800 (0.0168)	-0.4604	-0.2538	-0.1636	-0.1135
<b>Oldpeak</b>	-0.5403 (0.0115)	-0.6273	-0.3458	-0.2229	-0.1546
<b>Slope</b>	0.5793 (0.0977)	0.3570	0.1968	0.1269	0.0880
<b>Ca</b>	-0.7733 (0.0000)	-0.7908	-0.4360	-0.2811	-0.1949
<b>Thal</b>	-0.9004 (0.0019)	-0.5513	-0.3040	-0.1959	-0.1359

Table 2. Logistic regression coefficients (Mean/SD scaled data)

	Customary model (Pvalue)	Method 1	Method 2	Method 3	Method 4
<b>Intercept</b>	0.2319				
<b>Age</b>	-0.0419 (0.8323)	-0.0419	-0.0231	-0.0365	-0.0101
<b>Sex</b>	-0.8188 (0.0002)	-0.8172	-0.4505	-0.7106	-0.1966
<b>Cp</b>	1.0425 (0.0000)	1.0317	0.5688	0.8972	0.2483
<b>Trestbps</b>	-0.2409 (0.0596)	-0.2340	-0.1323	-0.2087	-0.0577
<b>Chol</b>	-0.2510 (0.2209)	-0.2297	-0.1266	-0.1997	-0.0553
<b>Fbs</b>	-0.0730 (0.9475)	-0.0755	-0.0416	-0.0657	-0.0182
<b>Restecg</b>	0.3668 (0.1806)	0.3711	0.2046	0.3228	0.0893
<b>Thalach</b>	0.3420 (0.0265)	0.3385	0.1866	0.2944	0.0815
<b>Exang</b>	-0.4276 (0.0168)	-0.4304	-0.2373	-0.3743	-0.1036
<b>Oldpeak</b>	-0.5950 (0.0115)	-0.6236	-0.3438	-0.5423	-0.1501
<b>Slope</b>	0.5568 (0.0977)	0.5641	0.3110	0.4905	0.1357
<b>Ca</b>	-0.7673 (0.0000)	-0.7983	-0.4402	-0.6943	-0.1921

Table 2 continued from previous page

	Customary model (Pvalue)	Method 1	Method 2	Method 3	Method 4
<b>Thal</b>	-0.5539 (0.0019)	-0.5676	-0.3129	-0.4936	-0.1366

To get  $b_2$ , (Equation 4) above result has to be divided by  $\pi\sqrt{3}$ , the numerical value of which is approximately 1.814. Hence, for Age the standardized coefficient becomes  $b_2 = (-0.004908) * (9.0821010) / 1.814 = -0.02457$ . To obtain the standardized coefficient by the third method (Equation 5) discussed in the previous section, the calculation is similar but instead of dividing by  $[\pi\sqrt{3}]$  the unstandardized coefficients are divided by  $[\pi\sqrt{3} + 1]$  which is equal to approximately 2.814. Therefore, for Age the calculation of the standardized coefficients is as follows:  $b_3 = (-0.004908) * (9.0821010) / 2.814 = -0.01584$ . The fully standardized fourth approach utilizes the value of the coefficient of determination ( $R^2$ ) to calculate the modified coefficients. This method, multiplies the first approach explained in equation 3 by  $R/S_{\logit(\hat{y})}$ . In this example, the value of the square root of  $R^2$  divided by the standard deviation of the  $\logit(\hat{y})$  was calculated to be 0.246434. So the modified coefficient for predictor Age changed in to  $b_4 = (-0.004908) * (9.0821010) * (0.246434) = -0.01098$ . Similar calculations have been done for all other variables and are presented in Table 1.

The column 'Customary model' in Table 1 refers to the calculated unstandardized coefficients from the logistic regression model. 'Method 1', 'Method 2', 'Method 3', and 'Method 4' represent the standardized coefficients computed using Equation 3, Equation 4, Equation 5, Equation 6 respectively. From the results in Table 1 it is evident that as the coefficients start from being partially standardized using method 1 to fully standardized in method 4, they seem to decrease in terms of magnitude. Techniques used in SAS have the closest values to the method suggested by Long. Predictor cp (chest pain) seems to have a comparatively higher relative effectiveness among the significant variables.

Table 3. Logistic regression coefficients (Median/MAD scaled data)

	Customary model (Pvalue)	Method 1	Method 2	Method 3	Method 4
<b>Intercept</b>	0.6920				
<b>Age</b>	0.0650 (0.9446)	0.0844	0.0465	0.0734	0.0177
<b>Sex</b>	-0.8415 (0.0007)	-0.8398	-0.4630	-0.7303	-0.1760
<b>Cp</b>	1.1343 (0.0000)	1.1226	0.6189	0.9762	0.2353
<b>Trestbps</b>	-0.1610 (0.5814)	-0.2322	-0.1280	-0.2019	-0.0487
<b>Chol</b>	0.1072 (0.9692)	0.1527	0.0842	0.1328	0.0320
<b>Fbs</b>	0.0673 (0.7820)	0.0696	0.0384	0.0606	0.0146
<b>Restecg</b>	0.1261 (0.7370)	0.1276	0.0703	0.1109	0.0267
<b>Thalach</b>	0.6672 (0.0004)	0.9032	0.4979	0.7854	0.1893
<b>Exang</b>	-0.7237 (0.0004)	-0.7284	-0.4016	-0.6335	-0.1527
<b>Oldpeak</b>	-0.8915 (0.0000)	-1.1391	-0.6280	-0.9906	-0.2387
<b>Slope</b>	0.9685 (0.0002)	0.9811	0.5409	0.8532	0.2056
<b>Ca</b>	-0.7982 (0.0000)	-0.8305	-0.4579	-0.7222	-0.1741
<b>Thal</b>	-0.6699 (0.0011)	-0.6864	-0.3784	-0.5969	-0.1439

In the next step, the target was to set up four different models using standardized coefficients calculated by these approaches and compare their performance based on prediction accuracy. To measure the prediction accuracy, the dataset was randomly divided into two sets; the testing set which contains 20% of the data and the training set which contains the rest of the data. The models were developed using the training set and the testing set was used to verify the overall accuracy. One of the major hurdles faced while setting up models to calculate their accuracies is that the predictors were measured using different scales. Hence, to make the comparison easier, the predictors were scaled before any kind of analysis was done. Firstly, all the variables were standardized using the mean and standard deviation of the corresponding independent variable. In addition, we computed standardized coefficients using the Median/ MAD standardized data.

Previously explained four methods of standardizing the coefficients were then applied to both of these scaled datasets. All

these calculations were done with the help of statistical software R. Outcomes of standardization of the coefficients are given in Table 2 and Table 3. Here in the Table 2 column 'Customary model' refers to the unstandardized coefficients of the dataset scaled by the mean and the standard deviation along with the four standardization methods for the coefficients in the following columns. Similarly, in the Table 3 column 'Customary model' refers to the unstandardized coefficients of the dataset scaled by the median and the mean absolute deviation along with the four standardization methods for the coefficients in the following columns.

### 2.3 Evaluation Criteria

To compare the performance of the models with different standardization techniques, we have used training accuracy and testing accuracy of the models. In predictive modeling for binary outcome variable the term accuracy refers to the fraction of correctly specified predictions made by the proposed model. The complete data is divided into two sets namely the training set and the testing set by a random split for instance in this analysis we have used 80% of the data for the training set and 20% for the test set. At first the prediction model is built on the training set and later applied on the test set to assess its prediction accuracy.

One predicament in this process is that, as the data are divided into training and testing sets randomly using R software, there is a chance of getting different results for different subsets which may result in bias. To solve this issue the complete process was repeated 1000 times and the average of these repetitions was taken for calculations of testing and training accuracy. Another criteria that is used in comparing the accuracy in binary predictive modeling is area under the Receiver operating characteristics (ROC). The plot represents the proportion of correctly specification events versus the proportion incorrect specification of the non-events for different probability cutoff's. A high area under the ROC curve indicates a better predictive accuracy.

### 2.4 Results

Table 4 shows the testing accuracy and training accuracy of the models constructed by applying each of the four coefficient standardization methods along with the model of unstandardized coefficients, which is represented by the 'Customary model' column.

Results indicate that the testing accuracy of the customary model was slightly higher than all standardized models for median/MAD scaled data. However, for the mean/SD scaled data, the testing accuracy for the customary model and the models for the 4 methods were similar. Similarly, the training accuracies were somewhat similar for the unstandardized and standardized coefficients. Moreover, method 4 was seen to have the lowest prediction accuracy among all four methods. On the other hand, by comparing the testing and training accuracies for mean/SD scaled data and median/MAD scaled data it can be seen that median/MAD scaled data has approximately 4% to 5% higher accuracy overall.

Table 4. Table for Testing and Training Accuracy

<b>Data</b>	<b>Customary model</b>	<b>Method 1</b>	<b>Method 2</b>	<b>Method 3</b>	<b>Method 4</b>
Mean					
Standardized (Test set)	0.8193	0.8218	0.8193	0.8221	0.8126
Median					
Standardized (Train set)	0.8754	0.8767	0.8646	0.8766	0.7813
Mean					
Standardized (Test set)	0.8576	0.8574	0.8540	0.8569	0.8472
Median					
Standardized (Train set)	0.9005	0.9000	0.8859	0.8987	0.8083

However, by taking a look at the AUCs for these models in Table 5 it can be seen that even though the unstandardized model had slightly different AUCs, there was no difference in AUCs of the models constructed from different standardization techniques. This indicates that in terms of distinguishing between the two diagnostic groups, all of these models show similar performance.

In terms of improving the sensitivity or specificity of the models the standardization techniques seem to have no significant effect. As the overall accuracy for the standardized models were lower than the un-standardized one, evidently the

sensitivity and specificity was also found to be less than the prior. Moreover, method 4 seems to have the higher sensitivity than all other models, which also means lower specificity than others.

Table 5. AUC's for Testing and training set

<b>Data</b>	<b>Customary model</b>	<b>Method 1</b>	<b>Method 2</b>	<b>Method 3</b>	<b>Method 4</b>
Mean					
Standardized (Test set)	0.8895	0.8899	0.8899	0.8899	0.8899
Median					
Standardized (Train set)	0.9230	0.9210	0.9210	0.9210	0.9210
Mean					
Standardized (Test set)	0.9262	0.9261	0.9261	0.9261	0.9261
Median					
Standardized (Train set)	0.9280	0.9279	0.9279	0.9279	0.9279

It is worth mentioning that the techniques used to scale the dataset seem to have some effect on improving the overall accuracy of the models. Test sets taken from the dataset for which the numerical variables were scaled using median/MAD standardization performed better than the one which was scaled using mean/standard deviation. For instance, for the customary model and the first three models, the testing accuracies were approximately 4% higher in the case of the dataset standardized by median/MAD Table 3. Additionally, from Table 5 it can be seen that the AUCs are slightly higher for the data which was standardized using median/MAD.

### 3. Discussion

The primary purpose of standardizing logistic regression coefficients is to set a ground on the basis of which the predictors can be ranked. The absolute value of the standardized coefficients enables one to order the independent variables in terms of importance. According to Menard (Menard S.W., 1995) standardized coefficients render a more precise idea than the un-standardized logistic regression coefficients. However, adapting such measures for the sake of interpretation may effect the overall performance of the model. In this study, the goal was to investigate how different standardization techniques effect the accuracy of the logistic regression model under study.

Different methods of standardizing the coefficients assist in explaining the variation in the dependent variable and allow one to compare their contributions. It was also investigated if standardizing the coefficients would change the performance of the model. From the results, it can be seen that if the standardized values are only used in the case of relative comparison of the predictors, there is not much difference between the four methods. The overall magnitude of the influence is comparatively lower for the 4<sup>th</sup> method but if the influences of the predictors were ranked, the ranking was found to be the same for all four methods.

By taking a closer look at the results it can be seen that standardizing the coefficients did not affect the overall prediction accuracy of the predictive logistic regression model. Similarly, no evidence was found that following a certain type of standardization technique would show better performance than the others; the unstandardized regression model, in general, had higher accuracy. However, method 4 would be a better approach compared to others, as method 2 suggested by Long and method 3 used in SAS, partially standardizes by only considering the predictors and does not include the outcome variable in calculation. Both of these methods make little difference to the outcomes thus are not recommended.

In essence, standardizing facilitates better interpretation and does not affect the predictive capacity of the model. This is evident from the AUC's computed for both unstandardized and standardized regression coefficients shown in Table 5. As the AUCs calculated from taking the average of multiple iterations, they turned out to be exactly equal for all standardization techniques, which was also similar to the un-standardized logistic regression model.

### 4. Conclusion

Logistic regression facilitates a wide range of techniques in conducting statistical analyses. In logistic regression like any other regression technique, the primary aim is to construct an equation based on the set of explanatory variables, which as a whole would explain the variation and predict the dependent variable better.

Therefore, it could be inferred that standardized coefficients can also be used for predictive modeling. Similarly, selecting

any specific method for standardizing the coefficients for interpretation is completely based on how one wants to interpret it. If the primary goal of conducting a logistic regression analysis is building up a predictive model which can also be used for comparing the predictor effects and does not affect the overall accuracy of the model, standardizing the regression coefficients may be advisable.

### Acknowledgment

The authors are grateful to the associate editor and the two referees for their valuable suggestions which have helped improve this paper.

### References

- Menard, S. W. (1995). *Applied logistic regression analysis*.
- Menard, S. (2000). Coefficients of determination for multiple logistic regression analysis. *The American Statistician*, 54(1), 17-24. <https://doi.org/10.2307/2685605>
- Greene, W. H. (1993). *Econometric Analysis*. (2nd & 4th ed.) Prentice Hall, New Jersey.
- Peng C. Y. J., Lee K. L., & Ingersoll G. M. (2002). An introduction to logistic regression analysis and reporting. *The journal of educational research*, 96(1), 3-14. <https://doi.org/10.1080/00220670209598786>
- Allison P. D. (1999). Comparing logit and probit coefficients across groups. *Sociological methods & research*, 28(2), 186-208. <https://doi.org/10.1177/0049124199028002003>
- Agresti A., & Kateri M. (2017). Ordinal probability effect measures for group comparisons in multinomial cumulative link models. *Biometrics*, 73(1), 214-219. <https://doi.org/10.1111/biom.12565>
- Agresti A., & Tarantola C. (2018). Simple ways to interpret effects in modeling ordinal categorical data. *Statistica Neerlandica*, 72(3), 210-223. <https://doi.org/10.1111/stan.12130>
- Abbott, R. D. (1985). Logistic regression in survival analysis. *American journal of epidemiology*, 121(3), 465-471. <https://doi.org/10.1093/oxfordjournals.aje.a114019>
- Long J. S. (1997). *Regression models for categorical and limited dependent variables*. Sage: vol. 7.
- Agresti, A. (2018). *An introduction to categorical data analysis*. John Wiley & Sons.
- Agresti, A., & Finlay, B. (1997). *Statistical Methods for the Social Sciences*.
- Kaufman, R. L. (1996). *Comparing effects in dichotomous logistic regression: A variety of standardized coefficients*. *Social Science Quarterly*. (p. 90-109).
- Menard S. (2004). Six approaches to calculating standardized logistic regression coefficients. *The American Statistician*, 58(3), 218-223. <https://doi.org/10.1198/000313004X946>
- Detrano R. (1989) *Cleveland heart disease database*.: VA Medical Center, Long Beach and Cleveland Clinic Foundation.
- Peterson, L. E., Nachemson, A. L., Bradford, D. S., Burwell, R. G., Duhaime, M., Edgar, M. A., ... & Willner, S. V. (1995). Prediction of progression of the curve in girls who have adolescent idiopathic scoliosis of moderate severity. Logistic regression analysis based on data from The Brace Study of the Scoliosis Research Society. *Journal of Bone and Joint Surgery-Series A*, 77(6), 823-827.
- Hassanipour, S., Ghaem, H., Arab-Zozani, M., Seif, M., Fararouei, M., Abdzadeh, E., ... & Paydar, S. (2019). Comparison of artificial neural network and logistic regression models for prediction of outcomes in trauma patients: A systematic review and meta-analysis. *Injury*, 50(2), 244-250.
- Kurt, I., Ture, M., & Kurum, A. T. (2008). Comparing performances of logistic regression, classification and regression tree, and neural networks for predicting coronary artery disease. *Expert systems with applications*, 34(1), 366-374.

### Copyrights

Copyright for this article is retained by the author(s), with first publication rights granted to the journal.

This is an open-access article distributed under the terms and conditions of the Creative Commons Attribution license (<http://creativecommons.org/licenses/by/4.0/>).